

# PATENT ABSTRACTS OF JAPAN

(11)Publication number : **05-081216**

(43)Date of publication of application : **02.04.1993**

(51)Int.Cl.

G06F 15/16

G06F 15/16

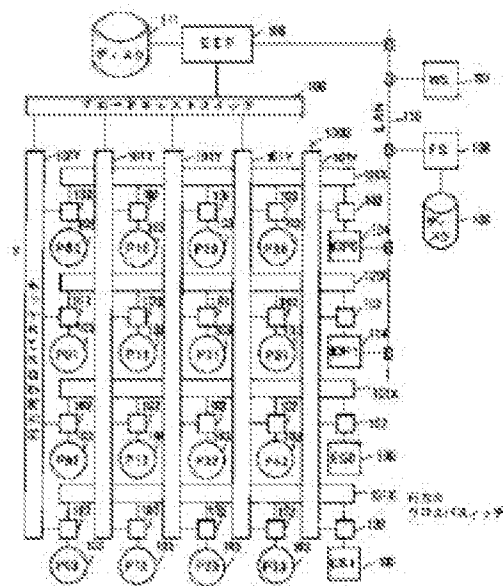
(21)Application number : **03-241092**

(71)Applicant : **HITACHI LTD**  
**HITACHI VLSI ENG CORP**

(22)Date of filing : **20.09.1991**

(72)Inventor : **NAKAKOSHI JUNJI**  
**HAMANAKA NAOKI**  
**CHIBA HIROYUKI**  
**HIGUCHI TATSUO**  
**SHUDO SHINICHI**  
**OGATA YASUHIRO**  
**TAKEUCHI SHIGEO**  
**TOBA TATSU**

## (54) PARALLEL PROCESSOR



### (57)Abstract:

**PURPOSE:** To provide a parallel processor capable of attaining a parallel processing managing function and an I/O function without reducing the number of processors for executing parallel processing and the loading method of the parallel processor.

**CONSTITUTION:** Cross bar switches 101 having  $2n+1$  ports are arrayed and a group of  $2n$  processors 103 is connected to the cross bar switches 101. Processors 104 to 106 for executing the parallel processing managing function and the I/O function are connected to the residual switches 101 out of respective switches 101 and respective processors are connected to the switches 101 through cross-over switches 102. Consequently parallel processing is executed by the  $2n+1$  processors and the parallel processing managing function and the I/O function are executed by other processors independently of the parallel processing. In the case of loading the parallel processor, each processor in a processor group connected to a cross bar switch in a certain dimension is

successively loaded.

(19)日本国特許庁(JP)

(12) 公開特許公報(A)

(11)特許出願公開番号

特開平5-81216

(43)公開日 平成5年(1993)4月2日

(51)Int.Cl.<sup>5</sup>

G 0 6 F 15/16

識別記号

4 0 0 Y 9190-5L

3 9 0 T 9190-5L

片内整理番号

F I

技術表示箇所

審査請求 未請求 請求項の数15(全 28 頁)

(21)出願番号 特願平3-241092

(22)出願日 平成3年(1991)9月20日

(71)出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(71)出願人 000233468

日立超エル・エス・アイ・エンジニアリング株式会社

東京都小平市上水本町5丁目20番1号

(72)発明者 中越 順二

東京都国分寺市東恋ヶ窪1丁目280番地  
株式会社日立製作所中央研究所内

(74)代理人 弁理士 小川 勝男

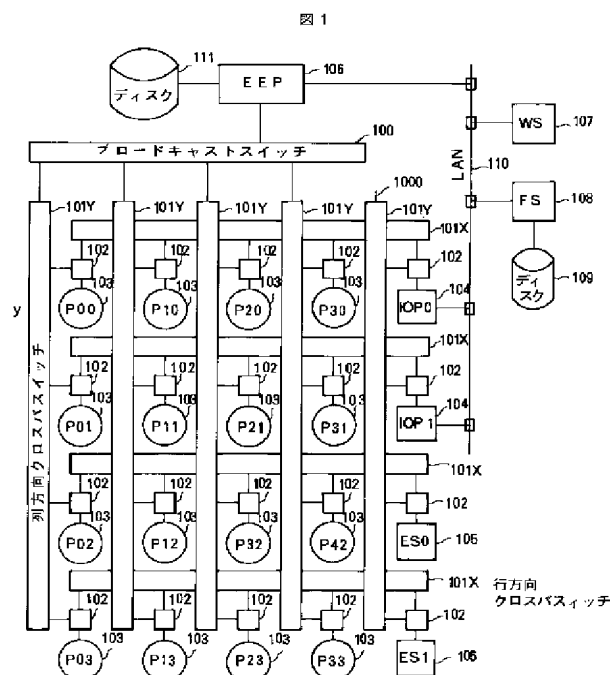
最終頁に続く

(54)【発明の名称】 並列プロセッサ

(57)【要約】 (修正有)

【目的】並列処理を実行するプロセッサの台数を減らすことなく、並列処理の管理機能や入出力機能を実現する並列プロセッサと、その実装方法を提供する。

【構成】2のべき乗+1個のポートを持つクロスバスイッチ101を設け、そのクロスバスイッチに2のべき乗個となるようにプロセッサ103群を配置する。また各クロスバスイッチ101の残りに並列処理の管理機能や入出力機能を実行するプロセッサ104、105、106を配置し、又各プロセッサとクロスバスイッチ101の接続には乗り換えスイッチ102を設ける。これにより並列処理の実行は2のべき乗個のプロセッサで実行し、並列処理の管理機能や入出力機能はそれ以外のプロセッサでその処理とは独立に実行できるようにした。並列プロセッサの実装では、ある1つの次元のクロスバスイッチ、そのクロスバスイッチに接続するプロセッサ群、および、そのプロセッサ群の中の1個のプロセッサに接続する。



## 【特許請求の範囲】

【請求項1】それぞれ並列処理プログラムを実行する複数の実行プロセッサと、それらを補助するための複数の補助プロセッサと、それらの実行プロセッサとそれらの補助プロセッサを接続するネットワークからなる並列プロセッサにおいて、

そのネットワークは、

それぞれ、並列にデータを転送可能な複数の部分ネットワークからなる複数の群の部分ネットワークと、

互いに異なる部分ネットワークに属する複数の部分ネットワークにそれぞれ含まれる複数の入出力ポートをそれぞれ接続する複数の乗り換えスイッチとを有し、

該複数の群の一つに属する部分ネットワークは、それぞれ2のべき乗プラス1個の入出力ポートを有し、

他の群の部分ネットワークの各々は、その群により定められた2のべき乗箇の入出力ポートを有し、

該複数の実行プロセッサは、該一つの群に属する複数の部分ネットワークの各々に含まれる2のべき乗個の入出力ポートの一つにそれぞれ接続された複数の乗り換えスイッチに接続され、

該複数の補助プロセッサは、該一つの群に属する複数の部分ネットワークの各々に含まれる該2のべき乗個の入出力ポート以外の一つの入出力ポートにそれぞれ接続された複数の乗り換えスイッチに接続されている並列プロセッサ。

【請求項2】各部分ネットワークは、クロスバスイッチからなる請求項1記載の並列プロセッサ。

【請求項3】該複数の補助プロセッサは、複数の入出力プロセッサを含む請求項1記載の並列プロセッサ。

【請求項4】該複数の補助プロセッサは、拡張記憶をそれぞれ有する複数の拡張記憶プロセッサを含む請求項1記載の並列プロセッサ。

【請求項5】該複数の補助プロセッサは、複数の入出力プロセッサと複数のフロントエンドプロセッサを含む請求項1記載の並列プロセッサ。

【請求項6】該複数の補助プロセッサは、少なくとも一つのフロントエンドプロセッサを含む請求項1記載の並列プロセッサ。

【請求項7】少なくとも1つのフロントエンドプロセッサと、

該他の群の内の一つの群に属する複数の部分ネットワークの各々に対応する入出力ポートおよび該フロントエンドプロセッサに接続された入出力ポートを有する他の一つの部分ネットワークとを更に有し、

該他の一つの群に属する複数の部分ネットワークの各々は、該他の部分ネットワークの一つの入力ポートに接続された一つの入出力ポートを更に有する請求項1記載の並列プロセッサ。

【請求項8】それぞれ並列処理プログラムを実行する複数の実行プロセッサと、それらを補助するための複数の

補助プロセッサと、それらの実行プロセッサとそれらの補助プロセッサを接続するネットワークからなる並列プロセッサにおいて、

そのネットワークは、

それぞれ、2の $i$ 乗（ $i$ は、正の整数）箇の第1の座標軸の座標点群の1つに対応して設けられ、それぞれ2の $j$ 乗（ $j$ は、正の整数）プラス1箇の入出力ポートを有する、2の $i$ 乗箇の第1種のクロスバスイッチと、

それぞれ、2の $j$ 乗箇プラス1箇の第2の座標軸の座標点群の1つに対応して設けられ、それぞれ少なくとも2の $i$ 乗箇の入出力ポートを有する、2の $j$ 乗箇プラス1箇の第2種のクロスバスイッチと、

それぞれ該第1の座標軸の座標点群と該第2の座標軸の座標点群との異なる組の1つに対応して配置され、それぞれ対応する組の第1の座標軸の座標点に対応する1つの第1種のクロスバスイッチ内の1つの入出力ポートと、該対応する組の第2の座標軸の座標点に対応する1つの第2種のクロスバスイッチ内の1つの入出力ポートとに接続され、それらの間でデータの転送をする複数の乗り換えスイッチを有し、

該複数のプロセッサは、該複数の乗り換えスイッチのうち、該第1の座標軸の座標点群と該第2の座標点群の内の連続する2の $j$ 乗箇の座標点との異なる組の1つに対応してそれぞれ配列された第2群の乗り換えスイッチの対応する1つにそれぞれ接続されている並列プロセッサ。

【請求項9】並列処理プログラムを並列に実行する、複数の実行プロセッサと、それらを補助するための複数の補助プロセッサと、それらの実行プロセッサとそれらの補助プロセッサを接続するネットワークからなる並列プロセッサにおいて、

そのネットワークは、

それぞれ、2の $i$ 乗（ $i$ は、正の整数）箇の第1の座標軸の座標点群と2の $k$ 乗（ $k$ は、正の整数）箇の第3の座標軸の座標点群との異なる組の1つにそれぞれ対応して設けられ、それぞれ2の $j$ 乗（ $j$ は、正の整数）プラス1箇の入出力ポートを有する、2の（ $i$ プラス $k$ ）乗箇の第1種のクロスバスイッチと、

それぞれ、該第2の座標軸の座標点群と該第3の座標軸の座標点群との異なる組のそれぞれに対応して設けられ、それぞれ少なくとも2の $i$ 乗箇の入出力ポートを有する、2の $j$ 乗箇プラス1掛ける2の $k$ 乗箇の第2種のクロスバスイッチと、

それぞれ、該第1の座標軸の座標点群と該第2の座標軸の座標点群との異なる組のそれぞれに対応して設けられ、それぞれ少なくとも2の $k$ 乗箇の入出力ポートを有する、2の $i$ 乗掛ける（2の $j$ 乗プラス1）箇の第3種のクロスバスイッチと、

それぞれ該第1の座標軸の座標点群と該第2の座標軸の座標点群と該第3の座標軸の座標点群との異なる組の1

つに対応して配置され、それぞれ対応する組の第1の座標軸の座標点に対応する1つの第1種のクロスバスイッチ内の1つの入出力ポートと、該対応する組の第2の座標軸の座標点に対応する1つの第2種のクロスバスイッチ内の1つの入出力ポートと、該対応する組の第3の座標軸の座標点に対応する1つの第3種のクロスバスイッチ内の1つの入出力ポートとに接続され、それらの間でデータの転送をする複数の乗り換えスイッチを有し、該複数のプロセッサは、該複数の乗り換えスイッチのうち、該第1の座標軸の座標点群と該第2の座標点群の内の連続する2のj乗箇の座標点と該第3の座標軸の座標点群との異なる組の1つに対応してそれぞれ配列された第1群の乗り換えスイッチの対応する1つにそれぞれ接続され、該複数の補助プロセッサの少なくとも一部の複数の補助プロセッサは、該複数の乗り換えスイッチのうち、該第1の座標軸の座標点群と該連続する複数の第2の座標軸の座標点以外の第2の座標点と該第3の座標軸の座標点群との異なる組の1つにそれぞれ対応する第2群の乗り換えスイッチの対応する1つに接続されている並列プロセッサ。

【請求項10】各第1種のクロスバスイッチと各第2種のクロスバスイッチは、互いに同じ回路からなる、請求項8または9記載の並列プロセッサ。

【請求項11】それぞれp箇の第1の座標軸の座標点群の1つに対応して設けられ、それぞれq箇の入出力ポートを有する、p箇の第1種のクロスバスイッチと、それぞれq箇の第2の座標軸の座標点群の1つに対応して設けられ、それぞれp箇の入出力ポートを有する、q箇の第2種のクロスバスイッチと、それぞれ第1から第3の入出力ポートを有し、それらの間でデータの転送をするp箇掛けるq箇の乗換えスイッチからなり、それぞれの乗換えスイッチは、該第1の座標軸の座標点群と該第2の座標軸の座標点群との異なる組の1つに対応して配列され、それぞれの第1の入出力ポートが対応する組の第1の座標軸の座標点群に対応する1つの第1種のクロスバスイッチ内の1つの入出力ポートに接続され、それぞれの第2の入出力ポートが対応する組の第2の座標軸の座標点群に対応する1つの第2種のクロスバスイッチ内の1つの入出力ポートに接続され、それぞれの第3の入出力ポートはデータ処理装置が接続され、該p箇の第1種のクロスバスイッチと該q箇の第2種のクロスバスイッチとp箇掛けるq箇の乗換えスイッチは、相互に接続された実質的に同一構造の複数の実装単位上に構成され、各実装単位は、1つの第1種のクロスバスイッチと、1つの第2種のクロスバスイッチと、それぞれ該1つの第1種のクロスバスイッチのp箇の入

出力ポートの対応する1つにその第1の入出力ポートに接続された、p箇の乗換えスイッチと、そのp箇の乗換えスイッチの内の1つの乗換えスイッチの第2の入出力ポートが上記1つの第2種のクロスバスイッチの1つの入出力ポートに接続され、該実装単位は、そのp箇の乗換えスイッチの内の残りのp箇の乗換えスイッチの第2の入出力ポートがそれぞれ他の対応する1つの実装単位上の1つの第2種のクロスバスイッチの1つの入出力ポートに接続されているクロスバネットワーク。

【請求項12】該pと該qはともに2のi乗プラス1（iは正の整数）である請求項1記載のクロスバネットワーク。

【請求項13】各実装単位上に該複数の装置の1部の装置がそれぞれ、その上の該p箇の乗換えスイッチのいずれか1つに含まれている第3の入出力ポートに接続されて、実装されている請求項1記載の並列プロセッサ。

【請求項14】それぞれ並列処理プログラムを実行する複数の実行プロセッサと、それらの実行プロセッサを接続するネットワークからなる並列プロセッサにおいて、少なくとも1つのフロントエンドプロセッサと、そのネットワークは、それぞれ、並列にデータを転送可能な複数の部分ネットワークからなる複数群の部分ネットワークと、互いに異なる部分ネットワークに属する複数の部分ネットワークにそれぞれ含まれる複数の入出力ポートをそれぞれ接続する複数の乗り換えスイッチとを有し、該複数の実行プロセッサは、該一つの群に属する複数の部分ネットワークの各々に含まれる2のべき乗個の入出力ポートの一つにそれぞれ接続された複数の乗り換えスイッチに接続され、該他の群の内の一つの群に属する複数の部分ネットワークの各々に対応する入出力ポートおよび該フロントエンドプロセッサに接続された入出力ポートを有する他の一つの部分ネットワークとを更に有し、該他の一つの群に属する複数の部分ネットワークの各々は、該他の部分ネットワークの一つの入力ポートに接続された一つの入出力ポートを更に有する並列プロセッサ。

【請求項15】各部分ネットワークは、クロスバスイッチからなる請求項14に記載の並列プロセッサ。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は並列プロセッサにおけるプロセッサの結合方式に関わり、特に並列処理の高速化に好適な並列プロセッサの構成と、その実装方法に関する。

【0002】

【従来の技術】従来、並列処理プログラムを実行する複数のプロセッサ（以下、実行プロセッサと呼ぶ）と、プ

ロセッサ間のデータ転送を可能にするネットワークにより構成される並列プロセッサは、たとえば特開昭63-124162号（以下、第1の従来例と呼ぶ）に示されている。この並列プロセッサにおいては縦2の $m$ 乗行と横2の $n$ 乗列（ $m$ 、 $n$ は正の整数）からなる2次元アレイ状に実行プロセッサを配置し、このアレイの各行（ $X$ 方向）、各列（ $Y$ 方向）ごとにクロスバスイッチを設ける。各実行プロセッサは $X/Y$ 方向のそれぞれの1つのクロスバスイッチに結合される。このため各実行プロセッサには1つの $X$ 方向のクロスバスイッチと1つの $Y$ 方向のクロスバスイッチを結合するために2つの入出力ポート（以下、単にポートと呼ぶ）を有している。以下、このような各次元ごとに設けたクロスバスイッチネットワークを部分クロスバネットワークと呼ぶ。

【0003】一方NCUBE社のNCUBE（以下、第2の従来例と呼ぶ）ではハイパキューブ構成のネットワークに入出力プロセッサを結合している（例えば、「32ビットマイクロプロセッサの全容 企業戦略・技術・市場動向」236、238ページ、日経マグロウヒル社発行、1986年12月）。これを実現するために各ノードプロセッサでは、ハイパキューブの次元数分に1つ加えた数のポートを持っている。ハイパキューブの次元数のポートには各ノードプロセッサを結合し、残りの1つには入出力プロセッサを結合する。即ち各ノードプロセッサは、補助プロセッサと接続するための専用のポートを有している。

【0004】また従来、部分クロスバネットワークを構成する並列プロセッサの実装については、例えばザ・セカンド・シンポジウム・オン・ザ・フロンティア・オブ・マシブリティ・パラレル・コンピューテーション（1988年10月）第463ページから第466ページ（The 2nd Symposium on the Frontiers of Massively Parallel Computation, October 10-12, 1988, pp463-466）、および電子情報通信学会技術研究報告、CPSY 89-45～58、1989年3月、第39ページから第44ページ（それぞれ第3、第4の従来例と呼ぶ）に示されている。

【0005】第3の従来例では2次元の部分クロスバネットワークで構成された並列プロセッサを、ローカルクロスバネットワークと呼ばれる1つのクロスバスイッチとそれに接続されるプロセッサ群を1つのウェハに実装する。またグローバルクロスバネットワークと呼ばれるもう1つのクロスバスイッチは前記ウェハとは別に、専用のLSIに実装する。このウェハと専用のLSIを複数個、組み合わせて並列プロセッサを構成する。

【0006】第4の従来例では3次元の部分クロスバネットワークで構成された並列プロセッサを、1つの面に対応する $X$ 方向と $Y$ 方向のクロスバスイッチ群とそれに接続されるプロセッサ群を1つの筐体に実装する。また $Z$ 方向のクロスバスイッチ群は専用の筐体に実装する。

これらの筐体を複数個、組み合わせて並列プロセッサを構成する。

【0007】

【発明が解決しようとする課題】各実行プロセッサを $X$ 方向と $Y$ 方向のクロスバスイッチに結合した上記第1の従来技術ではプロセッサの結合方式についてのみ開示されている。しかし並列処理を実行するためにはプロセッサ群への処理分配などの管理機能や入出力機能も不可欠であるが、これを行なう補助プロセッサについては言及していない。また、この補助プロセッサを上記プロセッサ群へ接続するためのネットワークを示していない。上記第2の従来例技術のように各実行プロセッサと補助プロセッサを直結したのでは、各実行プロセッサに新規に1つのポートを設け必要があり、かつ、そのポートでのメッセージ転送制御回路が複雑とならざるを得ない。また、この従来技術では補助プロセッサと各実行プロセッサを接続するネットワークを別に設けているが、この別のネットワークは実行プロセッサの台数の増大とともに複雑化せざるを得ない。

【0008】また第3、4の従来技術による並列プロセッサの実装においては、プロセッサの実装とは別に、クロスバスイッチ群を実装するための専用のLSI、筐体などが必要である。これによればプロセッサ台数が増えると、プロセッサとクロスバスイッチを接続するためのインタフェースが一ヶ所に集中し、そのインタフェースが多くなる。このためクロスバスイッチ群を実装する専用のLSI、筐体などの物理ピン数を越えることになり、そのインタフェースが接続できない場合がある。

【0009】本発明の第1の目的は補助プロセッサを実行プロセッサに接続する比較的構造の簡単なネットワークを有する並列プロセッサを提供することにある。

【0010】本発明の第2の目的はプロセッサとネットワークを接続するためのインタフェースを一ヶ所に集中させない並列プロセッサの実装方法を提供することにある。

【0011】

【課題を解決するための手段】上記第1の目的を達成するために、本発明では並列処理の管理機能や入出力機能を実行するための補助プロセッサを、任意の実行プロセッサ間でデータ転送を行なう部分クロスバネットワークに接続する。各次元のクロスバスイッチとして2のべき乗個のポートに1個のポートを追加したポート数（2のべき乗+1個）のポートを持つ部分クロスバスイッチを用いる。各次元が2のべき乗個となるように $n$ 次元直方体の実行プロセッサ群を配置し、それらのクロスバスイッチに乗換えスイッチを介して、また各クロスバスイッチの残りの1ポートのいずれかに並列処理の管理機能や入出力機能を実行する補助プロセッサに乗換えスイッチを介して接続する。

【0012】上記第2の目的を達成するために、本発明

では、ある1つの次元のクロスバスイッチ、そのクロスバスイッチに接続するプロセッサ群、および、そのプロセッサ群の中の1個のプロセッサに接続する、かつ、上記とは別の次元の全てのクロスバスイッチを、1つの実装単位とする。

#### 【0013】

【作用】並列処理は、従来通り、各次元が2のべき乗個となるように $n$ 次元立方体に配置した実行プロセッサ群で実行できる。並列処理を実行する実行プロセッサからは、プロセッサ間データ転送と同様にして補助プロセッサに対して通信が可能である。これにより各実行プロセッサにポートを追加することなく、また補助プロセッサを結合するためのネットワークの追加はごくわずかである。このため並列プロセッサの実装規模を小さくできる。

【0014】並列プロセッサの実装においては、クロスバスイッチ群を実装するための専用のLSI、筐体などは必要なく、この実装単位だけの組合せで可能である。これによりプロセッサとネットワークを接続するインタフェースを一ヶ所に集中させることなく実装できる。

#### 【0015】

【実施例】以下、本発明の一実施例を図を用いて説明する。図1に並列プロセッサのシステム構成を示す。本並列プロセッサにおいては行方向または列方向のクロスバスイッチ101（行方向クロスバスイッチを101x、列方向クロスバスイッチを101yと呼び、これらを区別しない場合には単に101と呼ぶ）、行方向クロスバスイッチ101xと列方向クロスバスイッチ101y間の乗り換えを行なう複数の乗り換えスイッチ102、 $i$ 行 $i$ 列に配列された並列処理を実行するプロセッサ（以下、計算クラスタと呼ぶ）103、入出力機能を実行する入出力プロセッサ（以下、入出力クラスタと呼ぶ）104、各計算クラスタ103で共有するデータや中間結果などを保持する拡張記憶プロセッサ（以下、ESクラスタと呼ぶ）105、並列処理の管理機能を実行するフロントエンドプロセッサ（以下、FEPと呼ぶ）106、一般的なユーザ環境を持つワークステーション（以下、WSと呼ぶ）107、ユーザファイルを管理するファイルサーバ（以下、FSと呼ぶ）108、ユーザファイルを保持するディスク装置109、WS107などのコンピュータ間の通信を高速に行なうローカルエリアネットワーク（以下、LANと呼ぶ）110により構成する。本実施例では、各クロスバスイッチ101は、2の $i$ 乗プラス1個のポートを有し、同じ構造を有し、さらに計算クラスタ103は2の $i$ 乗行2の $i$ 乗列に配列されるのが特徴である。

【0016】図1では一例として計算クラスタ103の台数を合計16台とし、それを $4 \times 4$ の2次元アレイ状に配置したものを示す。従って各クロスバスイッチは5個のポートを有する。しかし本発明は任意の値の $i$ に対

して有効である。

【0017】各計算クラスタ103は、行方向（以下、X方向と呼ぶ）と列方向（Y方向と呼ぶ）の各クロスバスイッチ101に乗り換えスイッチ102を介して接続する。図1では入出力クラスタ104とESクラスタ105はそれぞれ2台ずつあり、X方向のクロスバスイッチ101xの残りの1つのポートに接続する。FEP106は計算クラスタ103を接続するY方向の各クロスバスイッチ101yの残りの1つのポートにブロードキャストスイッチ100により接続される。ただし入出力クラスタ104とESクラスタ105とが接続されたY方向のクロスバスイッチ101y（図の右端のもの）の残りの1つのポート1000は、このクロスバスイッチ100には接続されない。これはブロードキャストスイッチ100のポートの数をクロスバスイッチ101と同じにしたため、このポート1000を接続するためのポートが存在しないからである。ここでは、このブロードキャストスイッチ100はクロスバスイッチ101と同じ構成である。各クロスバスイッチ101と各入出力クラスタ104、ESクラスタ105の接続にもそれらの間で乗り換えを行なう乗り換えスイッチ102を設ける。結局、本実施例では、2次元空間のX座標が0、1、2、3、Y座標が0、1、2、3からなる点群に対応して計算クラスタ103が配置される。各Y座標0～3の各点に対応して1つのX方向のクロスバスイッチが設けられ、各X座標に対応して1つのY方向のクロスバスイッチが設けられている。各入出力クラスタ104とESクラスタ105が $X=4$ 、 $Y=0 \sim 3$ の座標点群に対応して設けられていることになる。以下、これらのクロスバスイッチおよび乗り換えスイッチにより構成されるネットワークをクロスバネットワークと呼ぶ。

【0018】本並列プロセッサは、WS107などで作成された計算依頼をLAN110を介してFEP106で受け付ける。FEP106にWS107の機能も持ち、FEP106で直接、計算依頼を作成することも可能である。FEP106では、その計算依頼をスプーリングし、その計算依頼で指示されたプログラムの実行が可能なクラスタにそれを割当などの処理を行なう。また、これらのためにFEP106ではディスク装置111を保持する。各計算クラスタ103で実行するプログラムやデータにおいては、まずFEP106が各計算クラスタ103に対してプログラムをロードするための情報を配った後、各計算クラスタ103にプログラムの実行開始を指示する。そして各計算クラスタ103がそれを受け取ると、計算クラスタ103内のオペレーティングシステムにより、クロスバネットワークを介して、いずれかの入出力プロセッサ104に依頼し、FS108からその入出力プロセッサを介してプログラムやデータをロードする。ロードが終了すると、各計算クラスタ103はプログラムの実行を開始する。

【0019】各計算クラスタ103間でデータ転送が必要な時はクロスバネットワークを介して行なう。たとえば、計算クラスタP00からP33へのデータ転送では、まずP00から送出されたメッセージはP00に接続される乗り換えスイッチ102からX方向のクロスバスイッチ101x（図1の最上段の101x）へ送られる。そのX方向のクロスバスイッチ101xを通ったメッセージは、P30に接続される乗り換えスイッチ102によりY方向のクロスバスイッチ101y（図1の最右側から2番目の101y）に乗り換えられる。そして、Y方向のクロスバスイッチ101yからP33に接続される乗り換えスイッチ102に送られ、P33に届く。また同様に計算クラスタ103が計算の途中で入出力を必要とする時は入出力クラスタ104にクロスバネットワークを介して依頼する。依頼された入出力クラスタ104では依頼された内容に従い動作し、応答が必要な場合は、その結果を依頼元の計算クラスタ103に返す。またESクラスタ105では、各計算クラスタ103が起動される前にFEP106によりデータがロードされ、各計算クラスタ103は実行時にクロスバネットワークを介してそのデータをアクセスする。また計算クラスタ103が実行時に計算結果などを直接書き込んでよい。

【0020】図2に計算クラスタの構成を示す。計算クラスタ103は、一般的なプロセッサにクロスバネットワークとのインタフェースを持つように構成されていればよく、たとえば複数の計算処理を実行するプロセッシングユニット201、各計算クラスタ103間とのデータ転送を実行するネットワークインタフェースアダプタ202、プログラムやデータを記憶するメモリ203、計算クラスタ内のオペレーティングシステム（以下、OS）のブートデバイスやページングおよびスワップデバイスなどとして用いるディスク装置205および、そのディスク装置205との入出力インタフェースを持つ入出力アダプタ204により構成される。ここでは計算クラスタ103がマルチプロセッサ構成をとるようにプロセッシングユニット201が複数台あるが、1台であってもよい。

【0021】プロセッシングユニット201の処理において、計算クラスタ間のデータ転送要求、入出力要求、拡張記憶へのアクセス、FEPへの通信などがある場合は、ネットワークインタフェースアダプタ202へその処理を依頼する。ネットワークインタフェースアダプタ202ではこれらの要求を受け付けると、それに必要なメッセージを組立て、クロスバネットワークへ送り出す。またネットワークインタフェースアダプタ202ではクロスバネットワークから送られてくるメッセージを受け付け、それに対応した処理を行なう。その処理はたとえばプロセッシングユニット201にメッセージを渡すために割込みをかける、メッセージで指定された転送

データをメモリ203に書き込むなどである。

【0022】図3に入出力クラスタの構成を示す。入出力クラスタ104は、計算クラスタ103と同じようにネットワークインタフェースアダプタ202、メモリ203、入出力アダプタ204、ディスク装置205を持ち、さらに入出力を制御する入出力制御プロセッサ301、LAN110のインタフェースを持つLANアダプタ302により構成する。

【0023】入出力クラスタ104では、計算クラスタ103からの入出力要求メッセージをネットワークインタフェースアダプタ202で受け付けると、そのメッセージを入出力制御プロセッサ301に渡す。入出力制御プロセッサ301ではそのメッセージに従い、入出力アダプタ204、LANアダプタ302を介してディスク205に対して入出力動作を実行する。入出力動作後、計算クラスタ103への応答が必要な場合はネットワークインタフェースアダプタ202にその処理を依頼する。また、この処理以外にも入出力クラスタ104ではファイルサーバの機能も有する。

【0024】図4にESクラスタの構成を示す。ESクラスタ105はES制御プロセッサ401、ネットワークインタフェースアダプタ202、メモリ203、および拡張記憶（ES）402により構成される。

【0025】ES制御プロセッサ401は、後述するロック要求メッセージ、ES読出し要求メッセージなどの処理を行なう。また計算クラスタ103への応答がある場合は、ネットワークインタフェースアダプタ202へその処理を依頼する。ネットワークインタフェースアダプタ202では要求を受け付けると、それに必要なメッセージを組立て、クロスバネットワークへ送り出す。またネットワークインタフェースアダプタ202ではクロスバネットワークから送られてくるメッセージを受け付け、それに対応した処理を行なう。その処理はたとえばES制御プロセッサ401にメッセージを渡すために割込みをかける、メッセージで指定されたES書き込みデータをES402に書き込むなどである。またメモリ203はES制御プロセッサ401で実行するプログラムやデータなどを格納するものである。このES制御プロセッサ401は、一般的なマイクロプロセッサとメモリ制御ユニットで構成されてもよい。

【0026】次にメッセージの転送について説明する。図5に計算クラスタ間データ転送のメッセージフォーマットを示す。このメッセージは、メッセージを受け取るクラスタを表わす受信クラスタ番号501、ネットワークの転送方法を表わす転送制御502、メッセージの種別、受信クラスタ内での処理方法を表わすデータ転送コマンド503、受信クラスタに転送データ507が届いたことを知らせるためのデータ識別子アドレス504、受信クラスタに転送データ507が届いた時にその転送データ507をメモリ203に書き込むための転送デー

タ受信アドレス505、転送データ507の長さを示す転送データ長506および転送データ507により構成される。

【0027】各プロセッシングユニット201で実行されているプログラムがメッセージを転送する時には、プロセッシングユニット201はメッセージの各フィールドに設定する値をネットワークインタフェースアダプタ202に渡し、さらにメッセージ生成および送出の開始を指示する。ここで受信クラスタ番号501は図6に示すようにX方向/Y方向のクロスバスイッチのポートを決定するX系クロスバ番号とY系クロスバ番号で構成される。また各計算クラスタ間のデータ転送における送受信クラスタ番号は、プログラムあるいはシステムからは一般に順序付けられた番号にした方が、指定しやすい。たとえば16台の計算クラスタ構成の場合は0から15の番号が付けらる。しかし実際のハードウェアで使用するクラスタ番号は異なる。これは16台の計算クラスタ構成の場合、計算クラスタを指定するには4ビットであれば十分であるが、2のべき乗+1個のポートを持つクロスバスイッチを用いたことにより、ハードウェアとしてはX系クロスバ番号/Y系クロスバ番号にそれぞれ3ビットずつで合計6ビットが必要になる。これによると各計算クラスタは順序付けられた番号にならない。このため図12に示すような計算クラスタの番号変換を行なう必要がある。この変換は、たとえば計算クラスタ内のネットワークインタフェースアダプタ202で実現する。またFEP106、入出力クラスタ104の番号変換についても同図に示すように同様に実現される。

【0028】計算クラスタより送出されたメッセージは、クロスバスイッチ101および乗り換えスイッチ102がメッセージ内の受信クラスタ番号501および転送制御502を用いてルーティングを行ない、目的の受信クラスタに到着する。ここで転送制御502は図7に示すようにそのメッセージを全てのクラスタに転送するか否かを示すブロードキャスト制御とメッセージを転送する時にX方向あるいはY方向のどちらからを先に送出するかを決定する転送順序制御により構成される。この転送順序制御はクロスバネットワーク内のルーティングをたとえば、X方向からY方向に限る固定ルーティングとし、必ずその方向の転送だけにすれば必要ないが、各クラスタからFEP106へのデータ転送のように構成次第で、Y方向からX方向へメッセージを転送することが必要になるため転送順序制御を設けている。

【0029】メッセージが届いた受信クラスタでは、ネットワークインタフェースアダプタ202がそのメッセージのデータ転送コマンド503をみて、転送データ507を転送データ受信アドレス505に従いメモリ203に書き込む。また受信した転送データ507をメモリ203に書き込んだことを知らせるために、データ識別子アドレス504に従いメモリ203にその情報を書き

込む。受信クラスタ内のプロセッシングユニット201はデータ識別子アドレスによりメモリ203を監視し、転送データ507が届いたことを知る。

【0030】次にブロードキャストのメッセージ転送について説明する。各乗換えスイッチ102では対応する計算クラスタ103から転送されてきたメッセージ内のブロードキャスト制御がブロードキャストを示していることを知ると、まず、そのスイッチに接続されたX方向クロスバスイッチ101xにそのメッセージを転送する。上記X方向クロスバスイッチ101xでは、それに接続される全ての乗換えスイッチ102に対して、転送可能であることを確認してからそのメッセージを転送する。X方向クロスバスイッチ101xからブロードキャストのメッセージが転送されてきた乗換えスイッチ102の各々では、次にそのスイッチに接続されたY方向クロスバスイッチ101yにそのメッセージを転送する。以下同様にY方向クロスバスイッチ101yから乗換えスイッチ102にブロードキャストのメッセージを転送する。この乗換えスイッチ102では、このメッセージをそれに接続される計算クラスタ103に転送する。このようにしてブロードキャストのメッセージ転送を実行する。

【0031】次にFEP106と計算クラスタ103とのメッセージ転送について説明する。図1では、FEP106と各計算クラスタ103とのメッセージ転送を効率よく行うためにブロードキャストスイッチ100を1つ新たに設けている。これによりブロードキャストスイッチ100と一群のY方向のクロスバスイッチ101yを経由するだけで、FEP106といずれの計算クラスタ103との間でメッセージ転送が可能である。

【0032】たとえば図1のシステム構成における計算クラスタ103のP00からFEP106へのメッセージ転送では、計算クラスタ103間のメッセージ転送と同様にしてネットワークインタフェースアダプタ202からメッセージが送出される。このときのメッセージ内の受信クラスタ番号は、図12のクラスタ番号変換表で示すようにFEP106を示すX系クロスバ番号=4、Y系クロスバ番号=4に変換されている。またメッセージ内の転送順序制御は、図10に示すようにY方向クロスバスイッチからX方向クロスバスイッチへメッセージを転送するように、1に設定されている。計算クラスタ内のネットワークインタフェースアダプタ202から送出されたメッセージは、まずP00に接続される乗換えスイッチ102に送出される。その乗換えスイッチ102では、メッセージ内の転送順序制御により、それに接続されるY方向のクロスバスイッチ101yに転送する。そのY方向のクロスバスイッチ101yでは、図8で示すルーティング情報（詳細については後述する）に従い、メッセージ内の受信クラスタ番号のY系クロスバ番号に一致するポート番号804のポートにそのメッセ



ージを送出する。この場合、このポートにはブロードキャストスイッチ100が接続されている。ブロードキャストスイッチ100では、メッセージ内の受信クラスタ番号のX系クロスバ番号に一致するポート番号804のポートにそのメッセージを送出することにより、FEP106へのメッセージ転送が行なわれる。

【0033】FEP106から計算クラスタ103のP00へのメッセージ転送では、計算クラスタ103間のメッセージ転送と同様にしてネットワークインタフェースアダプタ202からメッセージが送出される。このときのメッセージ内の受信クラスタ番号は、図12のクラスタ番号変換表で示すように計算クラスタ103のP00を示すX系クロスバ番号=0、Y系クロスバ番号=0に変換されている。またメッセージ内の転送順序制御は、図10に示すようにX方向クロスバスイッチからY方向クロスバスイッチへメッセージを転送するように、0に設定されている。FEP106のネットワークインタフェースアダプタ202から送出されたメッセージは、まずブロードキャストスイッチ100に送出される。そのブロードキャストスイッチ100では、図8で示すルーティング情報（詳細については後述する）に従い、メッセージ内の受信クラスタ番号のX系クロスバ番号に一致するポート番号804のポートにそのメッセージを送出する。この場合、このポートにはY方向のクロスバスイッチ101yが接続されている。そのY方向のクロスバスイッチ101yでは、メッセージ内の受信クラスタ番号のY系クロスバ番号に一致するポート番号804のポートにそのメッセージを送出する。このポートには、P00に接続される乗換えスイッチ102が接続されており、これにより計算クラスタ103のP00へのメッセージ転送が行なわれる。

【0034】FEP106から全計算クラスタ103へのブロードキャスト転送では、計算クラスタ103間のメッセージ転送と同様にしてネットワークインタフェースアダプタ202からメッセージが送出される。このときのメッセージ内の受信クラスタ番号は何であってもよい。またメッセージ内の転送順序制御は、図10に示すようにX方向クロスバスイッチからY方向クロスバスイッチへメッセージを転送するように0に設定されている。さらにメッセージ内のブロードキャスト制御がブロードキャストを示している。FEP106のネットワークインタフェースアダプタ202から送出されたメッセージは、まずブロードキャストスイッチ100に送出される。そのブロードキャストスイッチ100では、ブロードキャストが指定されているため、全ポートに送出可能であることを確認して、そのメッセージを送出する。各々のY方向のクロスバスイッチ101yでは、ブロードキャストのメッセージが転送されてくると、それに接続される全ての乗換えスイッチ102に対して、メッセージが転送可能であることを確認して各乗換えスイッチ

102に転送する。

【0035】このようにFEP106から全計算クラスタ103へのブロードキャストでは、計算クラスタ間のそれとは異なり、X方向のクロスバスイッチ101xを介することなく、ブロードキャストスイッチ100と一群のY方向のクロスバスイッチ101yを経由するだけで、全ての計算クラスタ103にメッセージをブロードキャスト転送が可能である。また各計算クラスタ103からFEP106にメッセージを転送する場合にも、ブロードキャストスイッチ100を使うことによりX方向のクロスバスイッチ101xを介せずに転送できるので、ネットワーク内で実行されている他の計算クラスタ間のメッセージ転送の妨げを少なくできる。もしFEP106がたとえば図1の右上端の入出力プロセッサ104の位置にあった場合には、各計算クラスタ103から送出された複数のメッセージがそのFEP106が接続されるX方向のクロスバスイッチ101xに集中することになり、そのX方向のクロスバスイッチ101xを使用する他のメッセージ転送が妨げられるが、本実施例では、このような問題はない。

【0036】次にこのネットワークを構成するハードウェアの実現方法について説明する。図8にメッセージの転送に必要なルーティング情報を示す。このルーティング情報はシステムが立ち上がった時に予め設定される。各クラスタ103、104、105およびFEP106には各クラスタの自クラスタ番号801を持つ。この番号801は送信クラスタが応答を必要とする時など、メッセージにそれを追加する。受信クラスタではそれに対応するメッセージを送り返す時、その番号が受信クラスタ番号として使用される。また各乗換えスイッチ102にもそれに接続される各クラスタの自クラスタ番号802を持つ。この番号802はメッセージ内の受信クラスタ番号501と比較される。乗換えスイッチ102は、その結果によりクラスタやクロスバスイッチにメッセージを渡す。各クロスバスイッチ101には各クロスバスイッチ101がX方向/Y方向かの次元を示すクロスバスイッチ次元情報803および、各クロスバスイッチ101のそれぞれのポートがどこに接続されているかを示すポート番号804を持つ。クロスバスイッチ次元情報803はクロスバスイッチ101内のルーティングにおいて、メッセージ内の受信クラスタ番号501のX系クロスバ番号/Y系クロスバ番号のどちらを切り出すかを決定する。またポート番号804は受信クラスタ番号501から切り出されたX/Y系クロスバ番号と比較され、一致するポートへメッセージを送出するために使用される。

【0037】図9に乗換えスイッチの構成を示す。乗換えスイッチ102は、乗換えスイッチ102に接続されるクラスタ、クロスバスイッチから送出されるメッセージを一時保持するファーストイン・ファースト

ーアウト・メモリ（以下、FIFO）901、FIFO 901に保持されたメッセージを取り出し、そのメッセージの受信クラス番号501と転送制御502により転送先を決定するルーティング制御902および各ルーティング制御902からの要求に対して優先順位を決定し、各ポートにメッセージを送出するメッセージ送出制御903を各ポート対応に設け、それらと自クラス番号802により構成する。

【0038】FIFO901は読み出しと書き込み動作が非同期に行なえる。これには市販されているFIFO用LSIを用いてもよい。ルーティング制御902はFIFO読み出し制御921と転送先決定回路922により構成する。FIFO読み出し制御921では、FIFO901から出力されるEMPTY信号911が0であることを確認して、リードストローブ912を送出し、FIFO901に保持されたメッセージを読み出す。このEMPTY信号911はFIFO901に1個もデータが入ってない時は1で、それ以外は0である。

【0039】転送先決定回路922では、図10に示すように、読みだされたメッセージ内の受信クラス番号501のX系クロスバ番号、Y系クロスバ番号がそれぞれ各乗り換えスイッチ102に保持されている自クラス番号802のX系クロスバ番号、Y系クロスバ番号と比較される。この比較結果に基づきメッセージ内の転送制御502の転送順序制御の内容によりX方向から先に送出するか、あるいはY方向から先に送出するかを決定する。この転送先決定回路922の動作は、たとえば転送順序制御がX方向からY方向への転送が指定されている、即ち転送順序制御が0であると、X系クロスバ番号の比較結果が0、即ち不一致で、かつY系クロスバ番号の比較結果が0の場合X系クロスバスイッチに転送、X系クロスバ番号の比較結果が0で、かつY系クロスバ番号の比較結果が1、即ち一致の場合X系クロスバスイッチに転送、X系クロスバ番号の比較結果が1で、かつY系クロスバ番号の比較結果が0の場合Y系クロスバスイッチに転送、X系クロスバ番号の比較結果が1で、かつY系クロスバ番号の比較結果が1の場合クラスタに転送することを示す。そして転送先決定回路922では決定した転送先に、それに対応する信号923によりメッセージ送出制御903に送出することを指示する。

【0040】メッセージ送出制御903はプライオリティ制御931とセクタ932により構成される。プライオリティ制御931では、各ルーティング制御からの信号923を受け取り、それらの間で送出する優先順位を決め、それに対して送出を行なう。送出の制御は、メッセージ送出制御903が送出先に乗り換えスイッチと同様に設けられているFIFOから出力されるFULL信号904が0であることを確認して、ライトストローブ905を送出し、送出先のFIFOにメッセージを書き込む。このFULL信号904は送出先FIFOが一

杯で1個のデータも受け取れない時は1で、それ以外は0である。またプライオリティ制御931では一回の送出でメッセージの転送が終わらない場合、送り出すごとに信号933により対応するFIFO読み出し制御921に対してFIFO901から次のデータを読みだすことを指示する。

【0041】次にこの乗換えスイッチ102におけるブロードキャストのメッセージ転送の動作について説明する。ここではブロードキャストは、前述したように送出元計算クラスタ103、乗換えスイッチ102、X方向クロスバスイッチ101x、乗換えスイッチ102、Y方向クロスバスイッチ101y、乗換えスイッチ102、全計算クラスタ103の順番でメッセージの転送が行なわれる。

【0042】まず計算クラスタ103から送出されたメッセージは、乗換えスイッチ102に送出される。そして乗換えスイッチ102のルーティング制御902が、メッセージ内のブロードキャスト制御によりブロードキャストが指定されているかを判定する。ブロードキャストが指定されていると、ルーティング制御902は出力ポート番号レジスタ924に保持されているポート番号に従い、そのブロードキャストのメッセージをX方向のクロスバスイッチ101xが接続されるポートに出力する。

【0043】ここで出力ポート番号レジスタ924の内容は、ブロードキャストのメッセージをどこのポートに出力するかを示す情報である。この場合の出力ポート番号レジスタ924の内容は、計算クラスタからブロードキャストのメッセージを受け取るルーティング制御902はX方向クロスバスイッチ101xに出力するように、X方向クロスバスイッチ101xからブロードキャストのメッセージを受け取るルーティング制御902はY方向クロスバスイッチ101yに出力するように、Y方向クロスバスイッチ101yからブロードキャストのメッセージを受け取るルーティング制御902は計算クラスタに出力するように設定されている。また同様にX方向のクロスバスイッチ101xから送出されたブロードキャストのメッセージは、乗換えスイッチ102のルーティング制御902が判定を行ない、ブロードキャストが指定されているとルーティング制御902は出力ポート番号レジスタ924に保持されているポート番号に従い、Y方向のクロスバスイッチ101yが接続されるポートに出力する。また同様にY方向のクロスバスイッチ101yから送出されたブロードキャストのメッセージは、乗換えスイッチ102のルーティング制御902が判定を行ない、ブロードキャストが指定されているとルーティング制御902は出力ポート番号レジスタ924に保持されているポート番号に従い、計算クラスタが接続されるポートに出力する。ブロードキャストのメッセージ転送では、この動作以外は他の転送制御と同じで

ある。図11にクロスバスイッチの構成を示す。クロスバスイッチ101は、それに接続される乗り換えスイッチ102などから送出されるメッセージを一時保持し、それを要求に応じて取り出しを行なうメッセージ受信制御1101および、そのメッセージの受信クラスタ番号501と転送制御502により自ポートへ出力するかを決定し、各メッセージ受信制御1101からの要求の優先順位を決定し、各ポートにメッセージを送出するメッセージ送出制御1102を各ポート対応に設け、構成する。

【0044】メッセージ受信制御1101はメッセージを一時保持するFIFO901およびFIFO901に保持されたメッセージをレジスタ1111に取り出すFIFO読み出し制御1110により構成する。FIFO読み出し制御1110は、乗り換えスイッチのそれと同様にFIFO901から出力されるEMPTY信号1112が0であることを確認して、リードストロブ1113を送出し、FIFO901に保持されたメッセージを読みだし、レジスタ1111に保持する。このEMPTY信号1112はFIFO901に1個もデータが入っていない時は1で、それ以外は0である。そしてメッセージ受信制御1101は、FIFO901からメッセージをレジスタ1111に取り出したことを、信号1122によりメッセージ送出制御1102に知らせる。

【0045】メッセージ送出制御1102は、信号1122により各メッセージ受信制御1101から読みだされたメッセージ内の受信クラスタ番号501内のX系クロスバ番号/Y系クロスバ番号のどちらを切り出すかを、クロスバスイッチ次元情報803により決定する。そして、その受信クラスタ番号501から切り出されたX/Y系クロスバ番号とポート番号804とを一致回路1120で比較する。その比較結果が一致すれば、そのポートへメッセージを送出することを決定する。また一致しなければ何も動作はしない。このメッセージの送出では、送出が決定した各メッセージ受信制御1101間の送出順序がプライオリティ制御1121により決められる。またプライオリティ制御1121が送出先に乗り換えスイッチと同様に設けられているFIFOから出力されるFULL信号905が0であることを確認して、ライトストロブ905を送出し、送出先のFIFOにメッセージを書き込む。このFULL信号904は送出先FIFOが一杯で1個のデータも受け取れない時は1で、それ以外は0である。プライオリティ制御1121は送出先のFIFOにメッセージを書き込んだことを信号1123により、メッセージ受信制御1101に知らせる。メッセージ受信制御1101は、この信号1123により次のメッセージを読みだす。またプライオリティ制御1121では一回の送出でメッセージの転送が終わらない場合、送り出すごとに信号1124により対応するFIFO読み出し制御1110に対してFIFO9

01から次のデータを読みだすことを指示する。

【0046】次にこのクロスバスイッチ101におけるブロードキャストのメッセージ転送の動作について説明する。FIFO読み出し制御1110がメッセージ内のブロードキャスト制御を判定し、ブロードキャストのメッセージ転送であることを認識すると、FIFO901からメッセージをレジスタ1111に取り出したことを、信号1122により全てのメッセージ送出制御1102に知らせる。それぞれのメッセージ送出制御1102では、上記と同じ動作を行ない、送出先のFIFOにメッセージを書き込んだことを信号1123により、メッセージ受信制御1101に知らせる。メッセージ受信制御1101は、信号1123が全てのメッセージ送出制御1102から送られてくるのを判別する。信号1123が全てのメッセージ送出制御1102から送られてきたときに、次のメッセージを読みだす。全てから送られてこないときは、全ての信号1123が送られてくるまで次のメッセージを読みださない。またプライオリティ制御1121では一回の送出でメッセージの転送が終わらない場合、送り出すごとに信号1124により対応するFIFO読み出し制御1110に対してFIFO901から次のデータを読みだすことを指示する。

【0047】次にESクラスタのアクセスについて説明する。図13にESのアドレス空間を示す。ESは図1のようにクロスバスイッチに接続され、分散配置されている。各ソフトウェアからESは各計算クラスタのメモリと独立した1つのアドレス空間にみえる。各ESクラスタ105内のES402（図4）は、この1つの空間内の1部の連続領域が割り当てられている。各ESのデータは計算クラスタのメモリに転送してから使用する。計算クラスタから各ESクラスタのアクセスは計算クラスタ間データ転送と同様にクロスバネットワークを介する。各計算クラスタから指定されるESアドレスは、プログラムからESクラスタ間をまたがるアドレスとなる。このため実際のアクセスでは、このアドレスを複数ESクラスタのうち1個を指定するESクラスタ番号と、1個のESクラスタで保持できる最大容量のビット数からなるESクラスタ内アドレスに分けられる。このESクラスタ番号は計算クラスタから指定されるESアドレスの上位の数ビットから変換されて得られる。たとえばESアドレスが32ビット幅で、各ESクラスタ105内のES402（図4）の最大容量が256Mバイト（ $2^{28}$ ）であると、ESクラスタ番号に変換されるESアドレスのビット数は上位4ビットである。図14に、この場合における図1の構成時の変換方法を示す。この場合ESクラスタ105は2台だけがクロスバネットワークに接続されているため、ESアドレスの上位4ビットが（0000）<sub>2</sub>の場合ではESクラスタ番号のX系クロスバ番号=4、Y系クロスバ番号=2に、またESアドレスの上位4ビットが（0001）<sub>2</sub>の場合で

はESクラスタ番号のX系クロスバ番号=4、Y系クロスバ番号=3に変換される。この変換は、たとえば計算クラスタ103内のネットワークインタフェースアダプタ202(図2)で実現する。

【0048】このネットワークインタフェースアダプタ202のブロック図を20図に示す。これはプロセッシングユニット、メモリなどのインタフェース制御を行なうプロセッサバス制御回路2001、プロセッサバス制御回路2001からの指示でメッセージをクロスバネットワークに送出するメッセージ送出制御回路2002、およびクロスバネットワークから送られてきたメッセージを受信し、クラスタ内のプロセッシングユニット201に割込みをかけたり、メッセージで指定された転送データをメモリ203に書き込む処理などを制御するメッセージ受信制御回路2003により構成される。

【0049】さらにメッセージ送出回路2002は、メッセージの送出を制御するメッセージ送出制御2201、ハードウェアで生成するメッセージの転送データの最大バイト数を保持する最大転送データ長レジスタ2202(以下、MTUと呼ぶ)、プロセッシングユニットから指定される転送データ長を保持する転送データ長レジスタ2203(以下、LENと呼ぶ)、プロセッシングユニットから指定されるESアドレスを保持するESアドレスレジスタ2204(以下、ESAと呼ぶ)、プロセッシングユニットから指定される自クラスタ内のメモリをアクセスするためのメモリアドレスを保持する送信メモリアドレスレジスタ2205(以下、MS-Sと呼ぶ)、プロセッシングユニットから指定される受信側クラスタ内のメモリをアクセスするためのメモリアドレスを保持する受信メモリアドレスレジスタ2206(以下、MS-Rと呼ぶ)、プロセッシングユニットから指定される受信クラスタ番号を保持するクラスタ番号レジスタ2207(以下、CLNと呼ぶ)、MTU2202とLEN2203を比較する比較器2211、LEN2203からMTU2202を減算する減算器2212、ESA2204の上位ビットからESクラスタ番号に変換する変換回路2213、ESA2204とMTU2202を加算する加算器2214、MS-S2204とMTU2202を加算する加算器2215、MS-R2204とMTU2202を加算する加算器2216、CLN2207から受信クラスタ番号に変換する変換回路2217、および、メッセージ送出制御の指示により上記レジスタなどの内容を入力とし、メモリから転送データを読み出し、メッセージを組立て、クロスバネットワークに送出する送信回路2208から構成される。

【0050】この構成においてESにデータを書き込むための計算クラスタ103のメッセージの送出処理について述べる。ESを書き込むためにプロセッシングユニットからは、ネットワークインタフェースアダプタに対して、転送データ長がLEN2203に、ESにデータ

を書き込むためのESアドレスがESA2204に、ESに書き込むデータを自クラスタ内のメモリから読み出すためのメモリ主記憶アドレスがMS-S2205に、およびESに対するアクセス種別(書き込み要求)が設定される。ただしアクセス種別についてはメッセージ送出制御2201で保持される。またES書き込みに関係のないレジスタ(この場合ではMS-R2206、CLN2207である)には何も設定されない。

【0051】またMTU2202については、プロセッシングユニットから指定された転送データ長がとて長い場合、それをそのままメッセージとして組み立ててしまうと、クロスバネットワークでそのためのメッセージ転送でバスが長い時間、占有されてしまい、他のメッセージ転送に影響がでる。このためメッセージの最大転送データ長を予め決めておき、それよりメッセージが長くなる場合は、その最大転送データ長に従いメッセージを分割して送出する。この最大転送データ長を保持するのがMTU2202であり、この値はシステムが立ち上がる時に設定されている。

【0052】以下、メッセージ送出回路2002の動作手順を示す。

【0053】(1)ES書き込み処理に必要なレジスタが設定されたことをメッセージ送出制御2201が認識する。

【0054】(2)比較器2211においてプロセッシングユニットから指定された転送データ長LEN2203がハードウェアで生成するメッセージの最大転送データ長MTU2202以下であるかを判定し、その結果をメッセージ送出制御2201に渡す。

【0055】(3)メッセージ送出制御2201では、プロセッシングユニットから指定された転送データ長がハードウェアで生成するメッセージの最大転送データ長を越えている場合には、選択回路2209を制御し、転送データ長をMTU2203を選択し、送信回路2208に送る。またESA2204の上位のあるビット幅分は変換回路2213によりESクラスタ番号に変換され、送信回路2208に送られる。ESA2204の残りのビットおよびMS-S2205はそのまま送信回路2208に送られる。

【0056】(4)そしてメッセージ送出制御2201は、送信回路2208に対してメッセージの送信を行なうことを指示する。

【0057】(5)送信回路2208では、MS-S2206に従い、クラスタ内のメモリから転送データを転送データ長分を読み出し、それとメッセージ転送に必要な情報を組み合わせてメッセージとしてクロスバネットワークに送出する。

【0058】(6)送信回路2208がメッセージの送出を終えると、メッセージ送出制御2201に知らせる。

【0059】(7)メッセージ送出制御2201では送信回路2208からのメッセージの送出終了を知ると、ESA2204およびMS-S2205の値にMTU2202の値をそれぞれの加算器2214、2215にて加算を行なう。この加算結果はそれぞれのESA2204およびMS-S2205に設定する。

【0060】(8)この後、減算器2212にてLEN2203からMTU2202を減算し、その結果をLEN2203に設定する。

【0061】(9)(2)から(8)までをLEN2203で保持する値がMTU2202で保持する値以下になるまで繰り返す。

【0062】(10)LEN2203で保持する値がMTU2202で保持する値以下になると、メッセージ送出制御2201では選択回路2209を制御し、転送データ長をLEN2203を選択し、送信回路に送り、(2)から(8)までの処理を1回だけ行なうことにより、ESにデータを書き込むための計算クラスタ103のメッセージの送出処理を終える。

【0063】ここでESのアクセスが複数のESクラスタに渡る場合にも、たとえばESのアクセスを4Kバイト単位で、かつ4Kバイト境界とし、またハードウェアで生成するメッセージの最大転送データ長MTU2203を8Kバイトにすれば、このESのアドレス加算の結果を番号変換することにより、複数のESクラスタに渡っても、この構成で処理できる。

【0064】また、このネットワークインタフェースアダプタ202は、計算クラスタ103、入出力クラスタ104、ESクラスタ105、FEP106の中に構成され、計算クラスタ間のデータ転送でも、MTU2201、LEN2203、MS-S2204、MS-R2206、CLN2207などを用いてESのアクセスと同様に実現できる。

【0065】図15にESアクセスのメッセージフォーマットを示す。このメッセージには、大きく計算クラスタ間での排他制御を行なうためのロック処理、読み出し処理、および書き込み処理がある。これらのメッセージでは計算クラスタ間データ転送メッセージと同様に受信クラスタ番号、転送制御および、それぞれの処理に対応したデータ転送コマンドが付けられる。

【0066】ロック処理では、ES上にロック変数を定義し、それにより計算クラスタ間で排他制御を行なう。まず計算クラスタからのロック要求メッセージ(a)がESクラスタ内のネットワークインタフェースアダプタ202に届くと、そのメッセージをES制御プロセッサ401に渡す。ES制御プロセッサ401では、そのメッセージ内のESアドレスに従い、ES402からESデータを読み出す。そして、そのESデータとメッセージ内の比較データとを比較し、等しければメッセージ内の格納データをそのESアドレスに格納する。等しくな

ければ格納はしない。この比較、格納する処理を行なう間は、他のESアクセスに関する処理はネットワークインタフェースアダプタ202で行なわないようにする。またES制御プロセッサ401では、その比較結果に従い条件コードを作成する。そしてロック応答メッセージ(b)を計算クラスタに送り返すためにネットワークインタフェースアダプタ202に必要な情報を設定し、依頼する。ネットワークインタフェースアダプタ202では、ロック応答メッセージ(b)を組立て、ネットワークにメッセージを送出する。計算クラスタに受信されたロック応答メッセージ(b)は、その中のプロセッシングユニット201(図2)に渡され、その条件コードの内容によりロック処理が行なわれたことを知る。

【0067】またロックを解除するアンロック処理では計算クラスタからアンロック要求メッセージ(c)を送出する。ESクラスタでそのメッセージが届くと、ES制御プロセッサ401にそのメッセージを渡すことなく、ネットワークインタフェースアダプタ202がそのメッセージ内のESアドレスに従い格納データを書き込む。これによりアンロック処理は終了する。

【0068】ESの読み出し処理では、計算クラスタは前述したESアドレスの変換を行ない、読み出し要求メッセージ(d)を送出する。読み出し要求メッセージ(d)がいずれかのESクラスタ内のネットワークインタフェースアダプタ202(図4)に届くと、そのメッセージをES制御プロセッサ401に渡す。ES制御プロセッサ401ではそのメッセージを解釈し、ESデータの読み出しに必要な情報をネットワークインタフェースアダプタ202に設定し、ESデータ読み出し応答メッセージ(e)の送出を依頼する。ネットワークインタフェースアダプタ202では、ES制御プロセッサ401から設定された情報に従い、読み出しデータ長分だけESデータを読みだし、読み出し応答メッセージ(e)を作成し、ネットワークに送出する。読み出し応答メッセージ(e)を受け取った計算クラスタでは、ネットワークインタフェースアダプタ202(図2)がそのメッセージ内の読み出しデータ格納アドレスに従いメモリ203(図2)に書き込み、ES読み出しデータが届いたことをプロセッシングユニット301(図2)に知らせる。

【0069】ESの書き込み処理では、計算クラスタは前述したESアドレスの変換を行ない、書き込み要求メッセージ(f)を送出する。そのメッセージを受信したESクラスタ105内のネットワークインタフェースアダプタ202では、メッセージ内のESアドレスに従い、書き込みデータ長分だけES402に書き込みデータを書き込む。

【0070】次に入出力について説明する。入出力クラスタは図1のようにX方向クロスバスイッチ101xとY方向クロスバスイッチ101yに接続され、分散配置

されている。各入出力クラスタのアクセスは計算クラスタ間データ転送と同様にクロスバネットワークを介してメッセージ転送により行なわれる。入出力クラスタがX方向クロスバスイッチ101xとY方向クロスバスイッチ101yに接続されている理由は、全計算クラスタ103から各入出力クラスタ104に対してアクセスが可能にするためである。もし仮にX方向のクロスバスイッチに接続されている計算クラスタが同じクロスバスイッチに接続されている入出力クラスタだけにしかアクセスしないのであればY方向のクロスバスイッチはなくてもよい。

【0071】計算クラスタ103から入出力を行う場合は入出力命令により実行される。この入出力命令では、入出力クラスタ番号、入出力装置番号あるいはファイル名などが指定され、計算クラスネットワークインタフェースアダプタに渡される。このネットワークインタフェースアダプタでは、計算クラスタ間データ転送メッセージと同様に命令で指定された入出力クラスタ番号からハードウェアで使用する受信クラスタ番号に変換を行なう。またメッセージのハードウェアヘッダとして転送制御、転送データ長などを生成する。メッセージ内の転送データは、入出力装置に書き込むためのデータや入出力のプロトコルなどとして使用する情報であり、それとハードウェアヘッダをメッセージとしてクロスバネットワークに送出する。

【0072】メッセージを受信した入出力クラスタのネットワークインタフェースアダプタ202(図3)は入出力クラスタの入出力制御プロセッサ302にその制御を渡す。入出力制御プロセッサではメッセージの転送データの内容を解読し、入出力アダプタ204、LANアダプタ302などを介してメッセージの転送データの内容に対応して実際の入出力装置との間で処理を行なう。この処理についてはUNIXオペレーティングシステムなどの処理と同様でよい。また入出力動作後、計算クラスタ103への応答が必要な場合は、入出力制御プロセッサ302がネットワークインタフェースアダプタ202にその処理を依頼する。

【0073】次に本発明による並列プロセッサの実装について説明する。図16に図1の構成時の実装を示す。1601はまとめて実装する1つの単位であり、たとえば筐体、パッケージ、LSIなどのいずれかに対応する。この実装単位1601は1個のY方向のクロスバスイッチ101yとそれに接続されるそれぞれ4つの乗り換えスイッチ102と計算クラスタ103、および1個のX方向のクロスバスイッチ101xである。

【0074】また入出力クラスタ104とESクラスタ105の実装については、実装単位1601内の計算クラスタ103を入出力クラスタ104とESクラスタ105に置き換えて実装する。さらにFEP106とそれに接続する1個のブロードキャストスイッチ100は、

上記とは別に実装する。

【0075】この実装単位1601を4つと、その単位1601内の計算クラスタ103を入出力クラスタ104とESクラスタ105に置き換えた単位を1つと、FEP106とそれに接続する1個のブロードキャストスイッチ100の単位により、図1の構成の並列プロセッサが実現される。

【0076】この実装ではネットワークの構成要素であるクロスバスイッチ101を分散実装が可能である。これによりクロスバスイッチ群とクラスタ群間を接続するインタフェースが1ヶ所に集中することがない。またクロスバスイッチ群とクラスタ群に分けて、それぞれを集中実装した場合と比較し、乗り換えスイッチとクロスバスイッチとの接続するインタフェースにおいて、実装単位間を接続するインタフェースの本数を減らすことができる。たとえば実装単位1601が筐体であるならば筐体間を渡るケーブルの本数を減らすことができる。またクラスタの台数を拡張する場合にも、予めクロスバスイッチ101のポート数を増やしておけば、その実装単位1601で増設が可能である。さらにこの場合、実装単位1601内でクラスタとクロスバスイッチのインタフェースが閉じるため、その間のインタフェースは近距離になり、実装単位を渡るインタフェースに比べ、高速なクラスタ間データ転送が実現できる。この局所的に速いクラスタ間データ転送を考慮し、並列処理のクラスタへの割当てなどを行なえば、並列処理全体の性能を向上することも可能である。

【0077】また乗り換えスイッチ、各クラスタ、クロスバスイッチ間を接続するインタフェースは高速に実現可能な同期転送を用いることが考えられる。しかし、この同期転送ではクラスタの最大台数時のインタフェースによるディレイを考慮して設計を行なう必要があり、またクラスタの台数を拡張するごとにクロックスキュー合わせが必要になる。これを考えれば乗り換えスイッチ、各クラスタ、クロスバスイッチ間を接続するインタフェースは非同期転送にするのが好ましい。このため本実施例では前述したように、各クラスタと乗り換えスイッチ間、乗り換えスイッチとクロスバスイッチ間にFIFOを設け、非同期転送を実現する。

【0078】FIFOの動作は図17に示すように、書き込み動作ではFIFOから出力されるFULL信号が0であることを確認して、ライトストロブとFIFO書き込みデータを送出し、FIFOにそのデータを書き込む。このFULL信号は送出先FIFOが一杯で1個のデータも受け取れない時は1で、それ以外は0である。またFULL信号が1である時は、FULL信号が0になる、即ち読み出し動作が行なわれるまで書き込み動作を抑止する。一方、読み出し動作ではFIFOから出力されるEMPTY信号が0であることを確認して、リードストロブを送出し、FIFOに保持されたメッ

セージを読みだす。このEMPTY信号はFIFOに1個もデータが入ってない時は1で、それ以外は0である。またEMPTY信号が1である時は、EMPTY信号が0になる、即ち書き込み動作が行なわれるまで読み出し動作を抑止する。また、これらの書き込み動作と読み出し動作は上記の動作を守れば、非同期に行なうことができる。

【0079】これを用いれば各クラスタ、乗り換えスイッチ、クロスバススイッチ間をそれぞれ独立に動作することが可能になる。これにより、たとえば乗り換えスイッチにエラーが生じ、その解析のためサービスプロセッサ（以下、SVP）がその乗り換えスイッチに対して乗り換えスイッチ内のレジスタ群などをスキャンアウト動作をしている時も、それに接続される各クロスバススイッチやクラスタは、そのスキャンアウト動作を知る必要がなく、また、それにより停止する必要もなく通常動作を実行できる。これは、各クロスバススイッチやクラスタの動作が、乗り換えスイッチに対するデータ転送動作を非同期転送であるFIFOの書き込み／読み出し動作を用いており、そのFIFOから出力されるFULL/EMPTY信号により、抑止できるためである。このように乗り換えスイッチ、各クラスタ、クロスバススイッチ間を接続するインタフェースを非同期転送にすることにより、それらを分離して部分保守が可能になる。

【0080】以上、本実施例では16台の計算クラスタが2次元アレイ状に配置された場合について説明したが、各クラスタの台数、その次元数を限定する必要はなく同様な考え方で実施可能である。図18に3次元アレイ状に配置された並列プロセッサの構成のイメージを示す。ここで、1801は乗り換えスイッチ102と計算クラスタ103を、1802は乗り換えスイッチ102と入出力クラスタ104を、1803は乗り換えスイッチ102とESクラスタ105を、1つにしてそれぞれ図示している。この詳細な説明は省略するが、このように容易に実現できることがわかる。

【0081】またFEP106の接続においても新たなクロスバススイッチ100を設けるのではなく、図19に示すように図1の入出力クラスタ104、ESクラスタ105などの代わりに複数ある構成であってもよい。この構成ではFEP106の負荷が重く、1台で処理しきれない場合に、複数のFEP106を接続し、その負荷を分散するのに有効である。

【0082】なお以上の実施例においては、X方向のクロスバススイッチとY方向のクロスバススイッチとFEPに接続するクロスバネットワークを全て同一構造と仮定した。従って、それらのポートの数は全て同一とした。しかし本発明はこれらに限られるのではなく、X方向のクロスバススイッチとY方向のクロスバススイッチに接続するプロセッサの数は異なってもよく、従って、それらのポートの数は異なってもよい。またクロスバネット

ワークは、FEPから全あるいは各クラスタへの転送および各クラスタからFEPへの転送が行える必要最小限の機能があればよい。

#### 【0083】

【発明の効果】本発明によれば、2のべき乗+1個のポートを持つクロスバススイッチを設け、そのクロスバススイッチに2のべき乗個となるように実行プロセッサ群を配置し、クロスバススイッチの残りの1ポートに並列処理の管理機能や入出力機能を実行する補助プロセッサを配置することができるので、各実行プロセッサにポートの追加、補助プロセッサを結合するためのクロスバススイッチの追加は必要なくなり、並列プロセッサの実装規模を小さくできる。さらに並列処理の実行は並列処理の管理機能や入出力機能とは独立に実行でき、並列処理を実行するプロセッサの台数を減らすことなく2のべき乗台とし、並列プロセッサを構成できる。

【0084】また本発明では、ある1つの次元のクロスバススイッチ、そのクロスバススイッチに接続するプロセッサ群、および、そのプロセッサ群の中の1個のプロセッサに接続する全ての、上記とは別の次元のクロスバススイッチを、1つの実装単位とするため、並列プロセッサの実装においては、クロスバススイッチ群を実装するための専用のLSI、筐体などは必要なく、この実装単位だけの組合せで可能である。これによりプロセッサとネットワークを接続するインタフェースを一ヶ所に集中させることなく実装できる。

#### 【図面の簡単な説明】

【図1】本発明による2次元構成時の並列プロセッサの一実施例を示すシステム構成図。

【図2】図1の構成に用いる計算クラスタの概略構成図。

【図3】図1の構成に用いる入出力クラスタの概略構成図。

【図4】図1の構成に用いる拡張記憶クラスタの概略構成図。

【図5】本発明による計算クラスタ間データ転送メッセージフォーマット図。

【図6】本発明によるデータ転送メッセージ内の受信クラスタ番号を示す図。

【図7】本発明によるデータ転送メッセージ内の転送制御を示す図。

【図8】図1の構成に用いるデータ転送のルーティング情報を示す図。

【図9】図1の構成に用いる乗り換えスイッチの概略構成図。

【図10】図9の構成に用いる転送先決定回路を示す図。

【図11】図1の構成に用いるクロスバススイッチの概略構成図。

【図12】図1の構成に用いる計算クラスタ番号変換を

示す図。

【図13】本発明による拡張記憶の空間を示す図。

【図14】図1の構成に用いる拡張記憶クラス番号変換を示す図。

【図15】本発明による拡張記憶アクセスのメッセージフォーマット図。

【図16】図1の構成における並列プロセッサの実装を示す図。

【図17】本発明で用いるFIFOのタイムチャート図。

【図18】本発明による3次元構成時の並列プロセッサの一実施例を示すシステム構成イメージ図。

【図19】本発明によるフロントエンドを複数結合した2次元構成時の並列プロセッサの一実施例を示すシステム構成イメージ図。

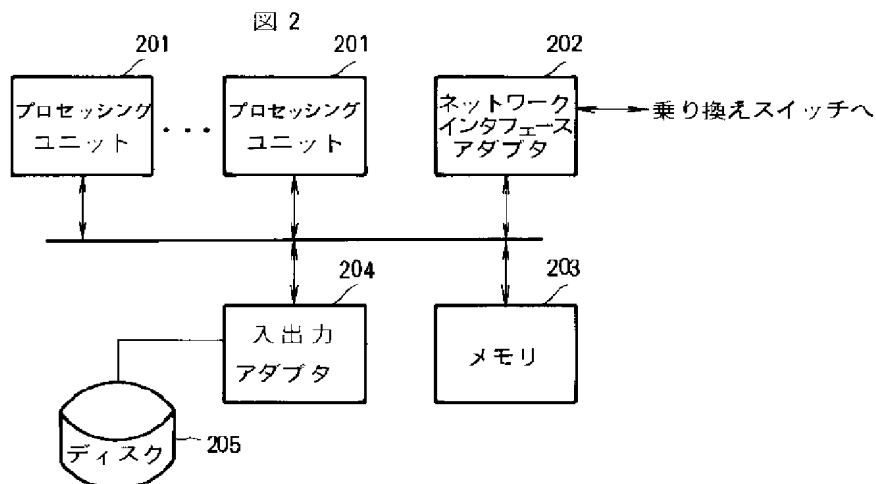
【図20】本発明によるネットワークインタフェースア

ダプタのブロック図。

【符号の説明】

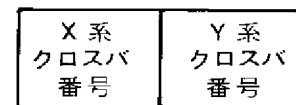
101…クロスバススイッチ、102…乗り換えスイッチ、103…計算クラスタ、104…入出力クラスタ、105…拡張記憶クラスタ、106…フロントエンド、107…ワークステーション、108…ファイルサーバ、109…ディスク装置、110…ローカルエリアネットワーク（LAN）、111…ディスク装置、201…プロセッシングユニット、202…ネットワークインタフェースアダプタ、203…メモリ、204…入出力アダプタ、205…ディスク装置、301…入出力制御ユニット、302…LANアダプタ、901…ファーストインファーストアウトメモリ（FIFO）、904…FIFOから出力されるFULL信号、911…FIFOから出力されるEMPTY信号、1601…まとめて実装する1つの単位。

【図2】



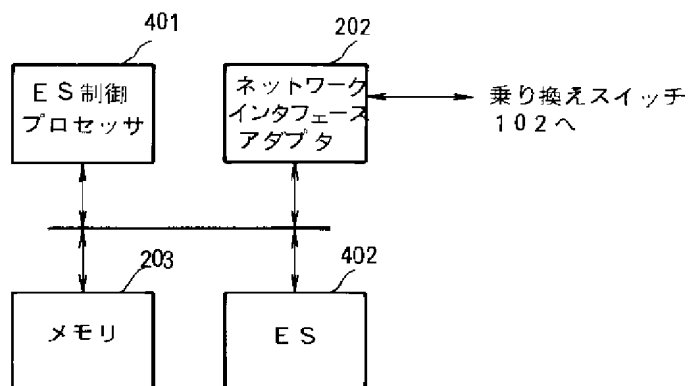
【図6】

図 6



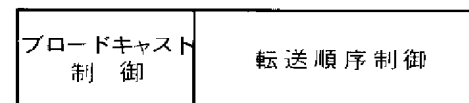
【図4】

図 4



【図7】

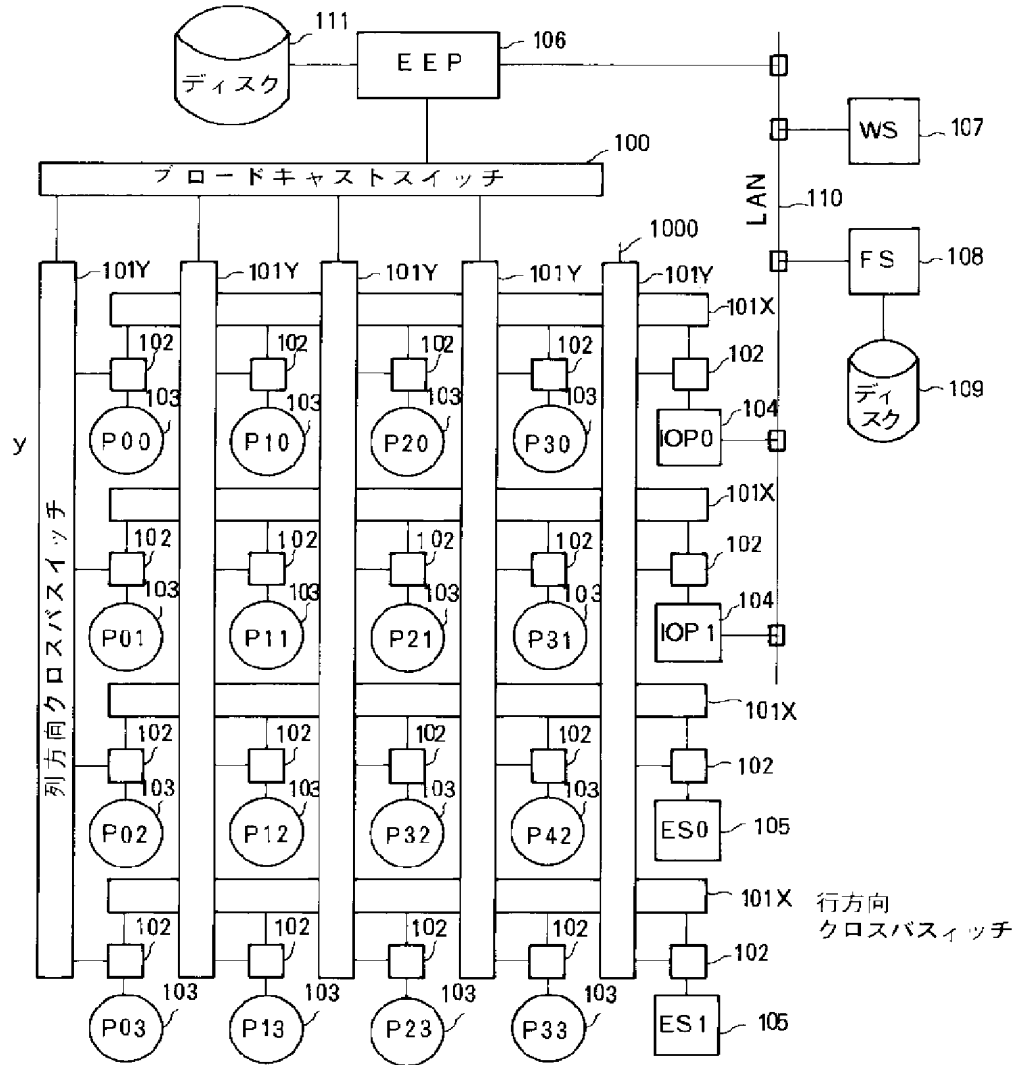
図 7





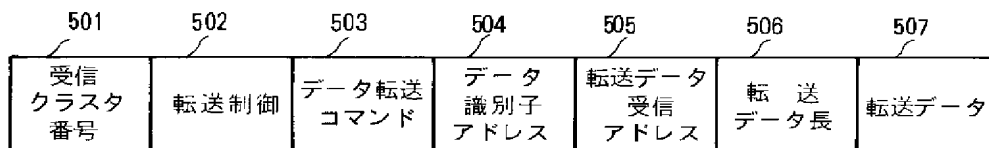
【図1】

図 1



【図5】

図 5



【図3】

【図14】

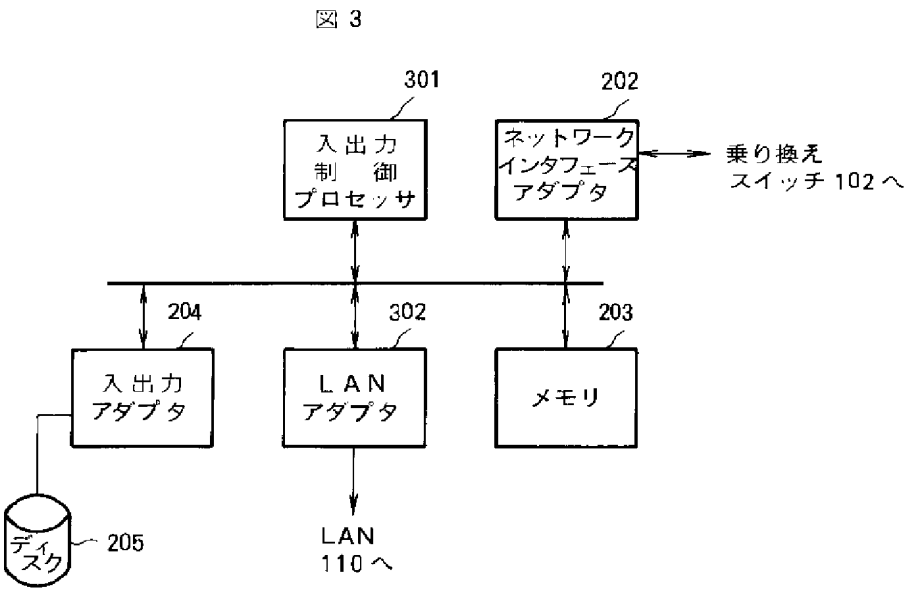


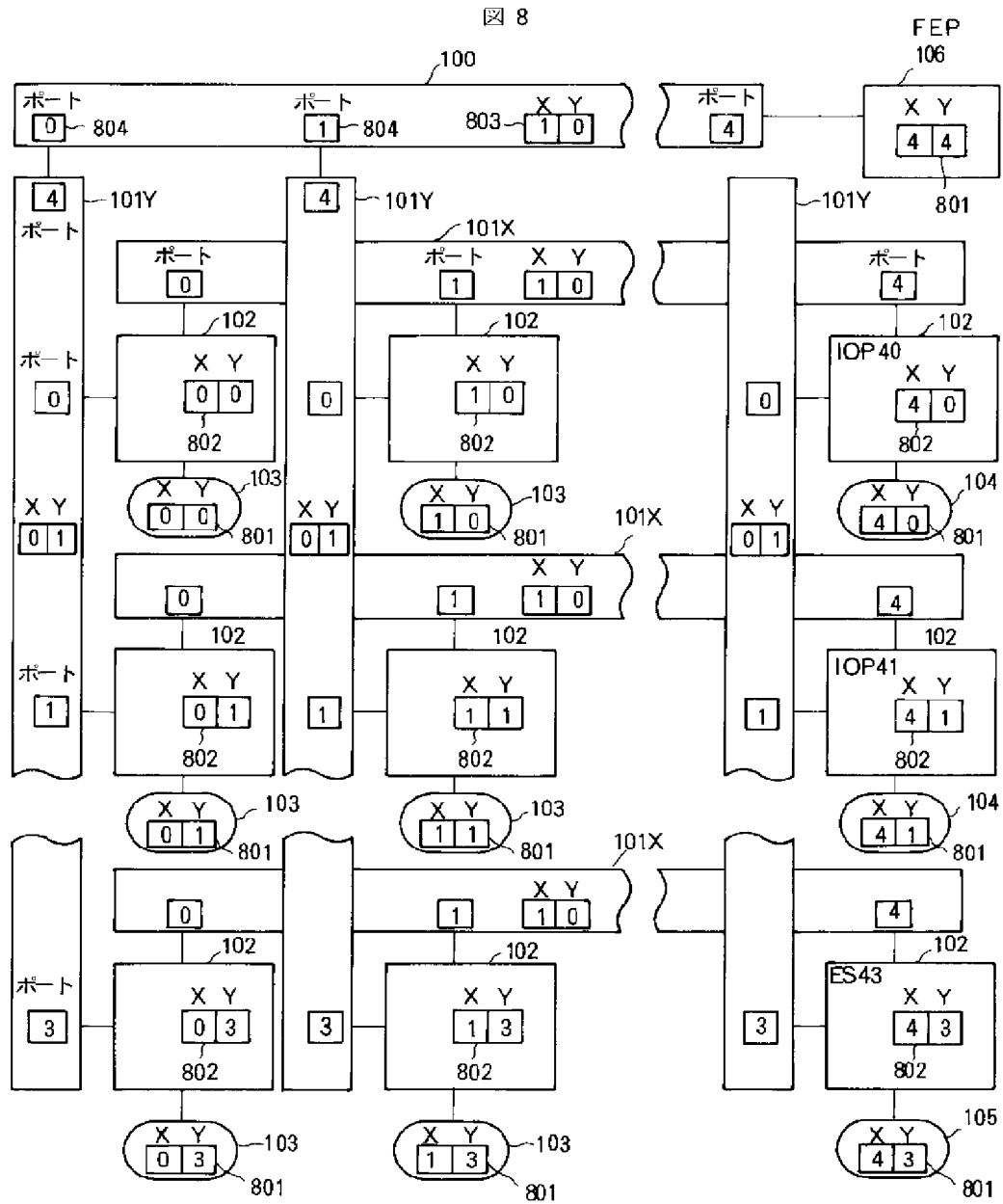
図 14

入 力	出 力	
	X	Y
0 0 0 0	4	2
0 0 0 1	4	3

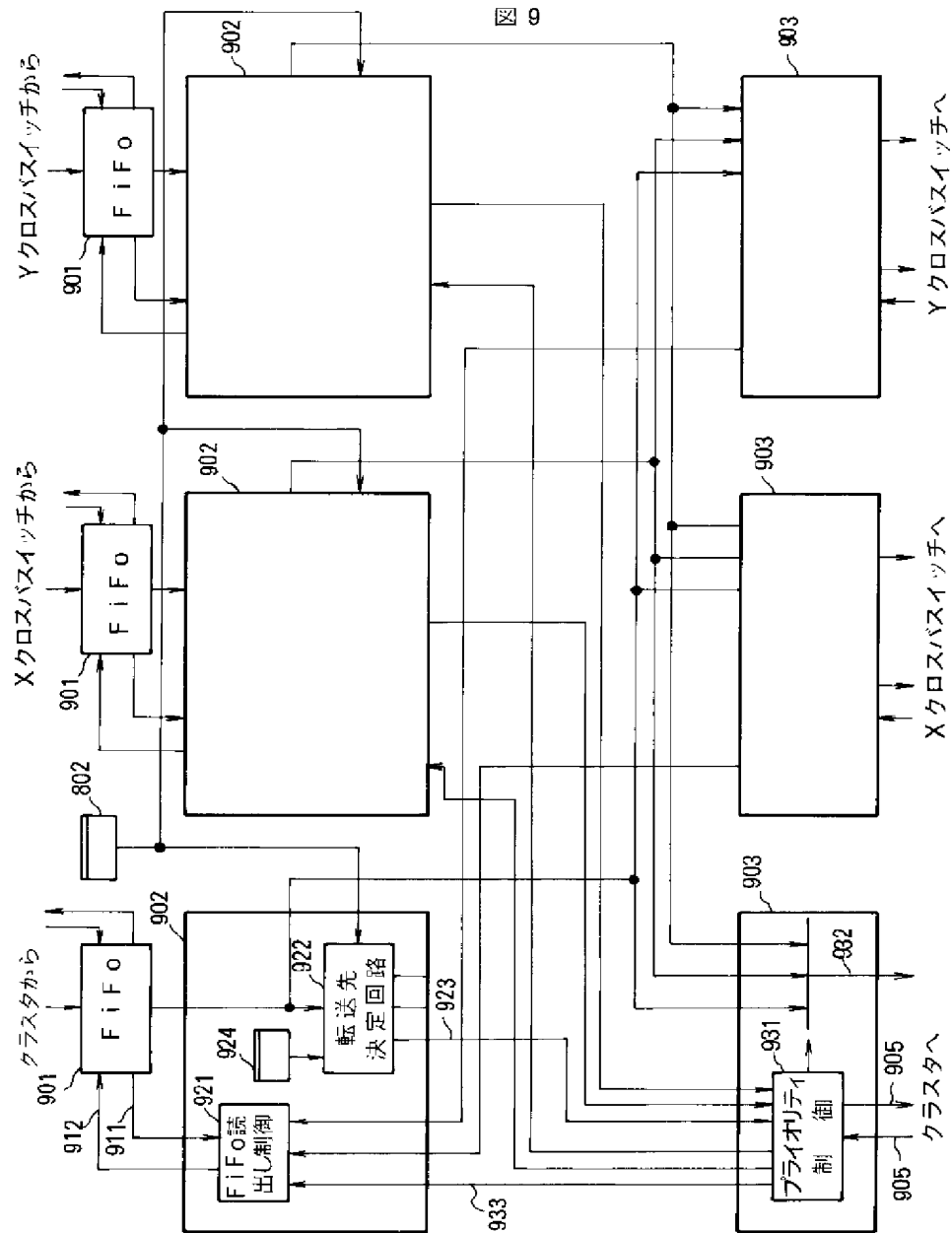
【図10】

受信クラス番号の比較結果		転 送 順 序 制 御	
		0 ( X系クロスバから転送 )	1 ( Y系クロスバから転送 )
X 系	Y 系		
0 ( 不一致 )	0 ( 不一致 )	X系クロスバに転送	Y系クロスバに転送
0	1 ( 一致 )	同 上	同 上
1 ( 一致 )	0	Y系クロスバに転送	X系クロスバに転送
1	1	クラスタに転送	クラスタに転送

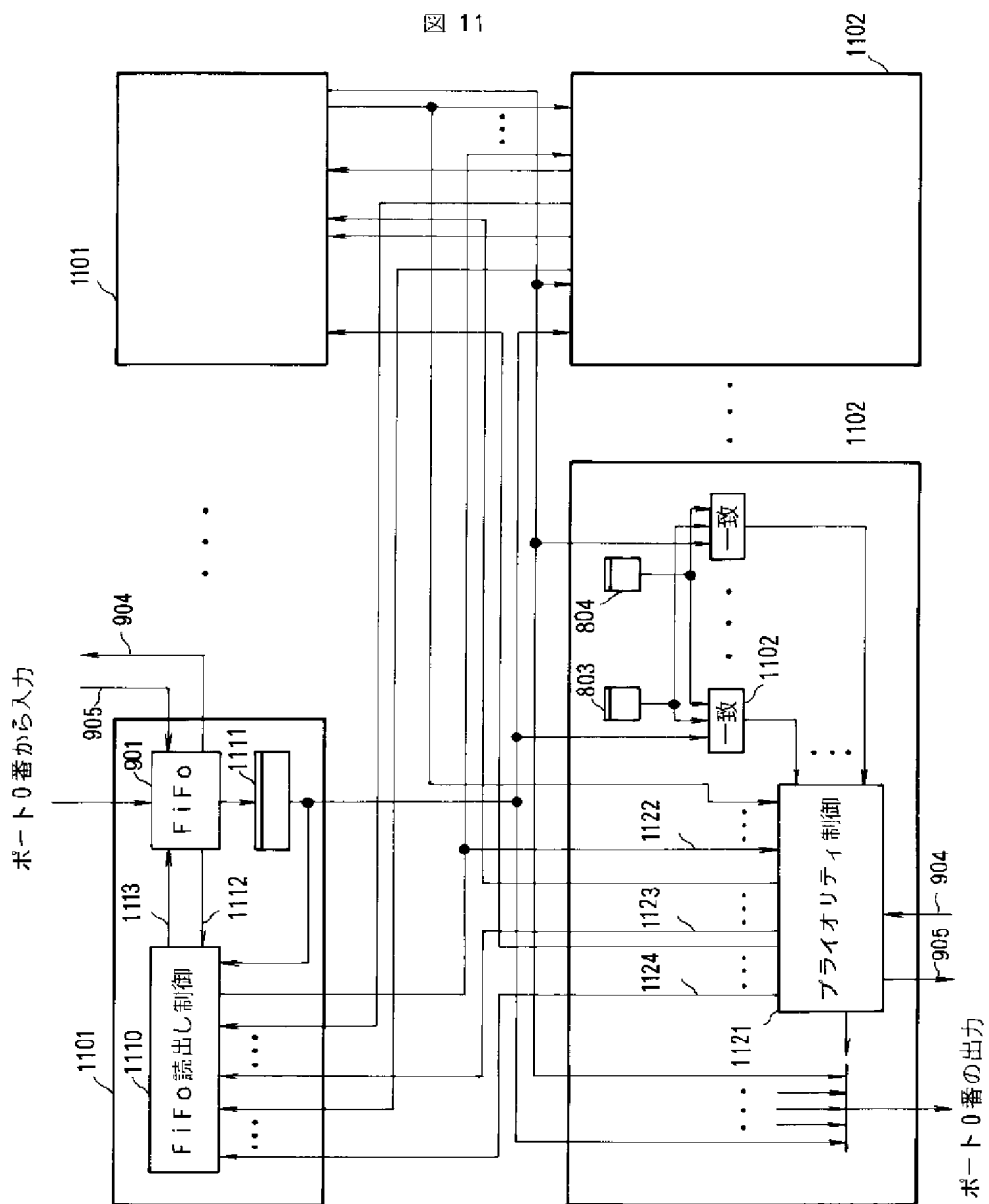
【図8】



【図9】



【例 1 1】



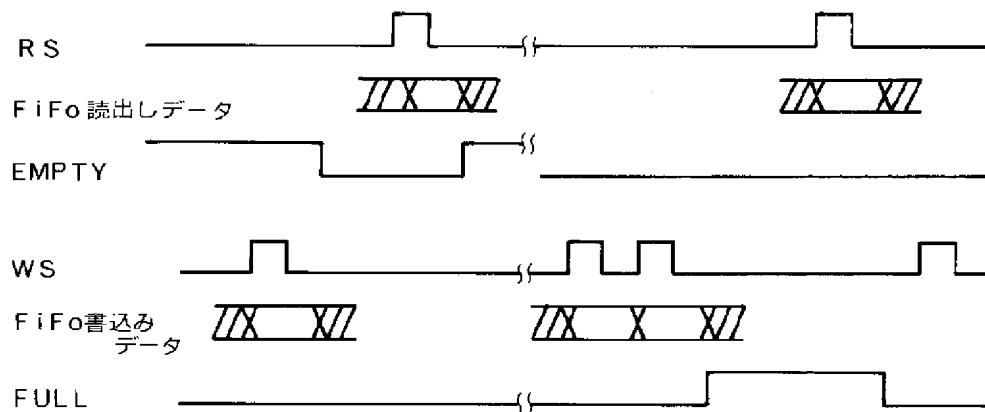
【図12】

図 12

入 力		出 力	
意 味	ビットパターン	X	Y
P 0 0	0 0 0 0 0	0 0 0	0 0 0
P 1 0	0 0 0 0 1	0 0 1	0 0 0
P 2 0	0 0 0 1 0	0 1 0	0 0 0
P 3 0	0 0 0 1 1	0 1 1	0 0 0
P 0 1	0 0 1 0 0	0 0 0	0 0 1
P 1 1	0 0 1 0 1	0 0 1	0 0 1
P 2 1	0 0 1 1 0	0 1 0	0 0 1
P 3 1	0 0 1 1 1	0 1 1	0 0 1
P 0 2	0 1 0 0 0	0 0 0	0 1 0
P 1 2	0 1 0 0 1	0 0 1	0 1 0
P 2 2	0 1 0 1 0	0 1 0	0 1 0
P 3 2	0 1 0 1 1	0 1 1	0 1 0
P 0 3	0 1 1 0 0	0 0 0	0 1 1
P 1 3	0 1 1 0 1	0 0 1	0 1 1
P 2 3	0 1 1 1 0	0 1 0	0 1 1
P 3 3	0 1 1 1 1	0 1 1	0 1 1
IOP0	1 0 0 0 0	1 0 0	0 0 0
IOP1	1 0 0 0 1	1 0 0	0 0 1
F E S	1 0 0 1 0	1 0 0	1 0 0

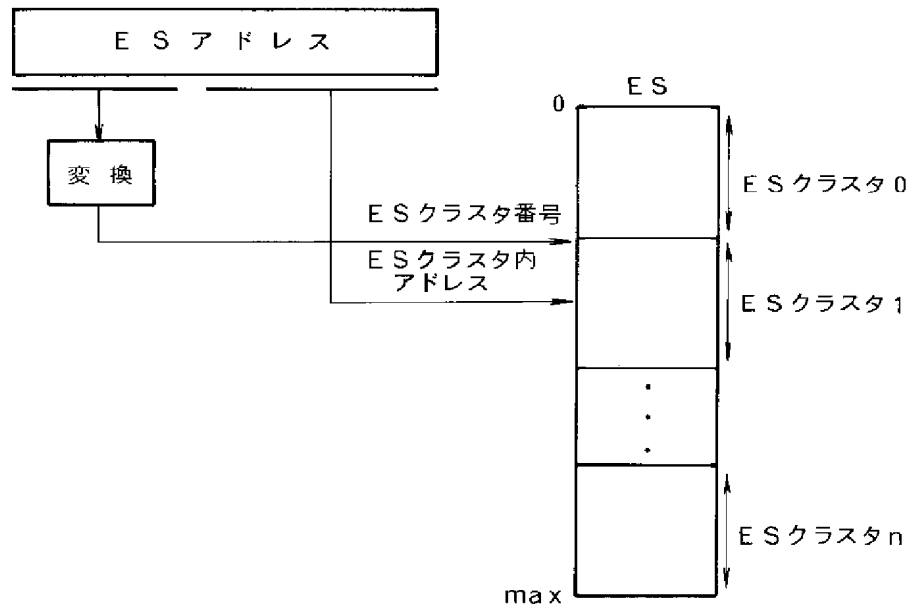
【図17】

図 17



【図13】

図 13



【図15】

図 15

## (a) ロック要求

受信(ES) クラスタ番号	転送制御	ロック要求	送信 クラスタ番号	ESアドレス	比較データ	格納データ
------------------	------	-------	--------------	--------	-------	-------

## (b) ロック応答

受信 クラスタ番号	転送制御	ロック応答	条件コード
--------------	------	-------	-------

## (c) アンロック要求

受信(ES) クラスタ番号	転送制御	アンロック 要 求	ESアドレス	格納データ
------------------	------	--------------	--------	-------

## (d) 読出し要求

受信(ES) クラスタ番号	転送制御	読出し要求	送信 クラスタ番号	ESアドレス	読出しデータ 格納アドレス	読出し データ長
------------------	------	-------	--------------	--------	------------------	-------------

## (e) 読出し応答

受信 クラスタ番号	転送制御	読出し応答	読出しデータ 格納アドレス	読出し データ長	読出しデータ
--------------	------	-------	------------------	-------------	--------

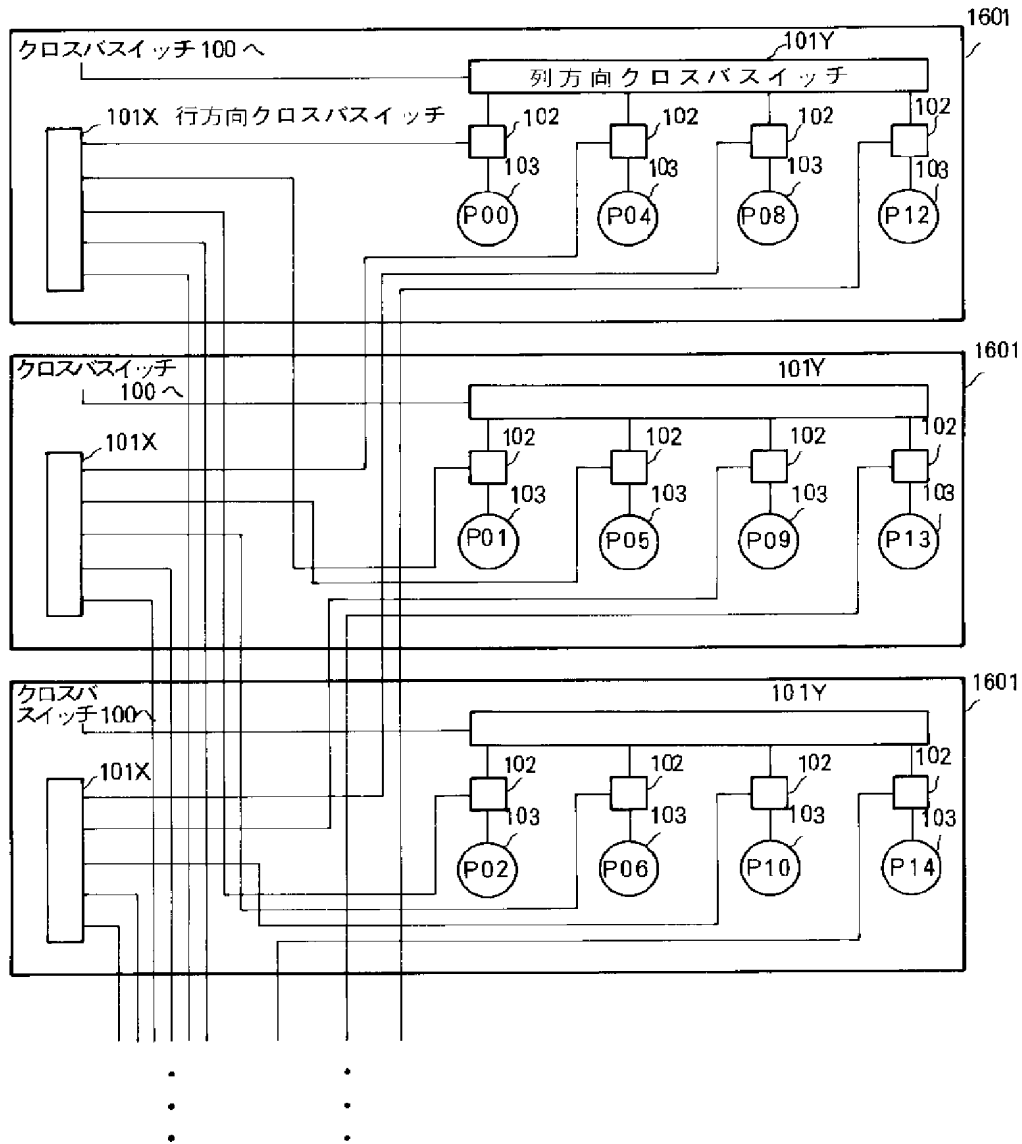
## (f) 書込み要求

受信(ES) クラスタ番号	転送制御	書込み要求	ESアドレス	書込み データ長	書込みデータ
------------------	------	-------	--------	-------------	--------



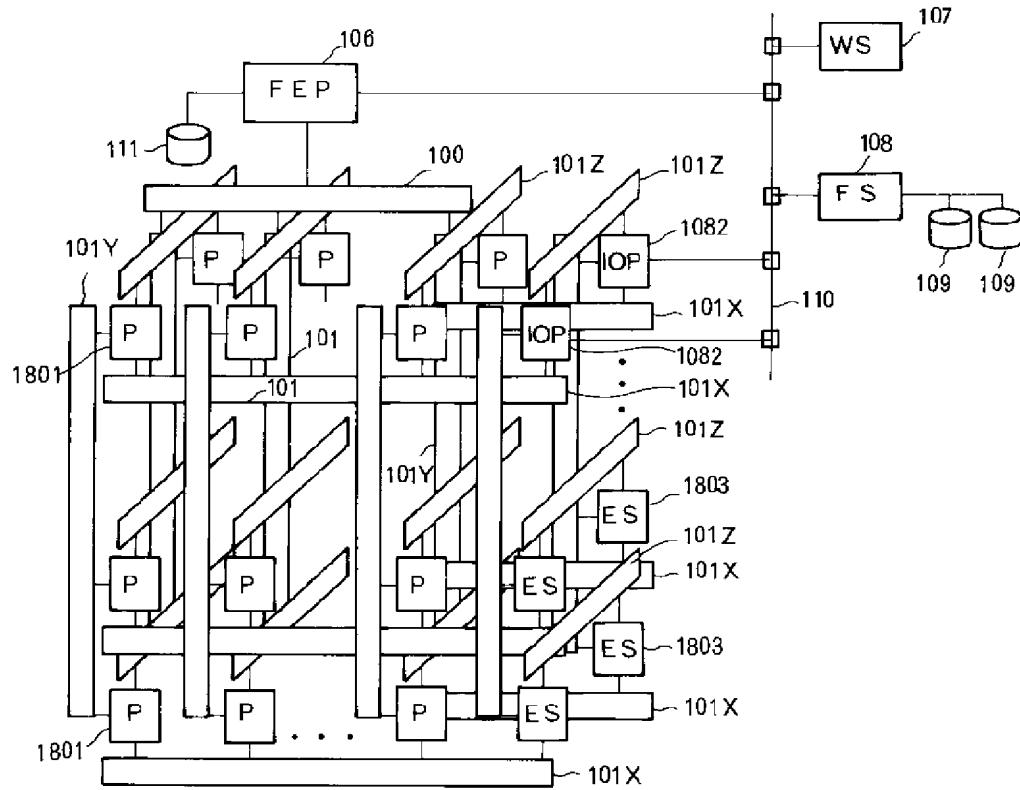
【図16】

図 16



【图 18】

☒ 18



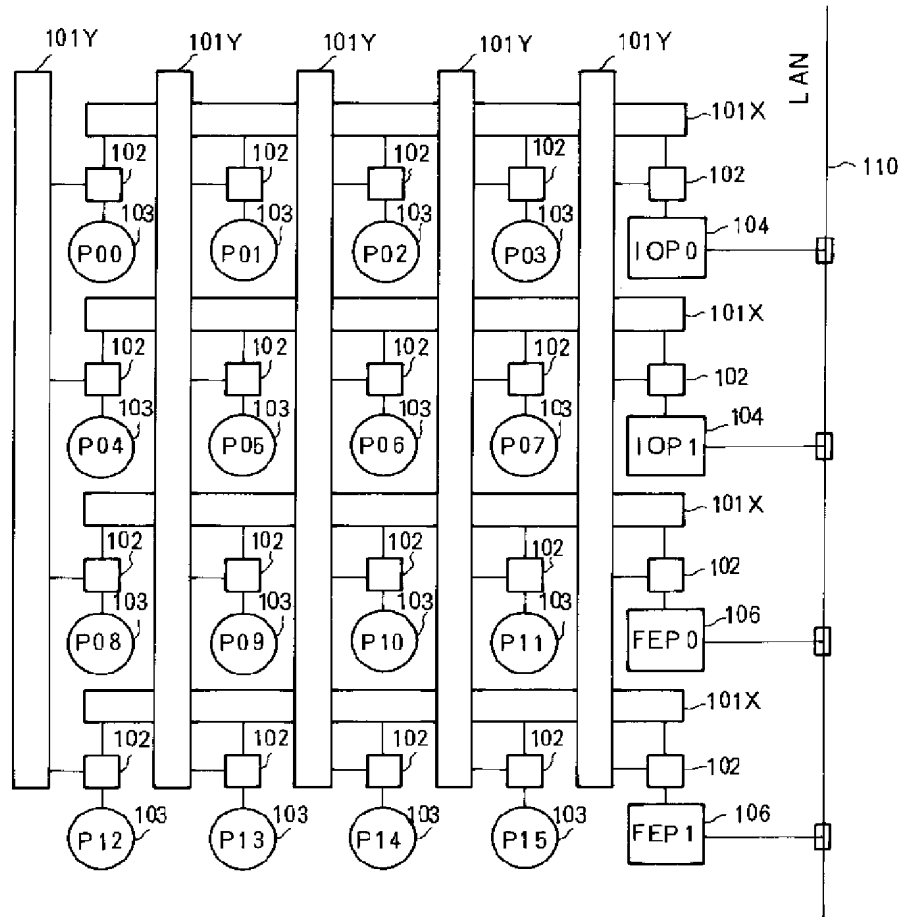
1801: 計算クラスタと乗り換えスイッチ

1802：入出力クラスタと乗り換えスイッチ

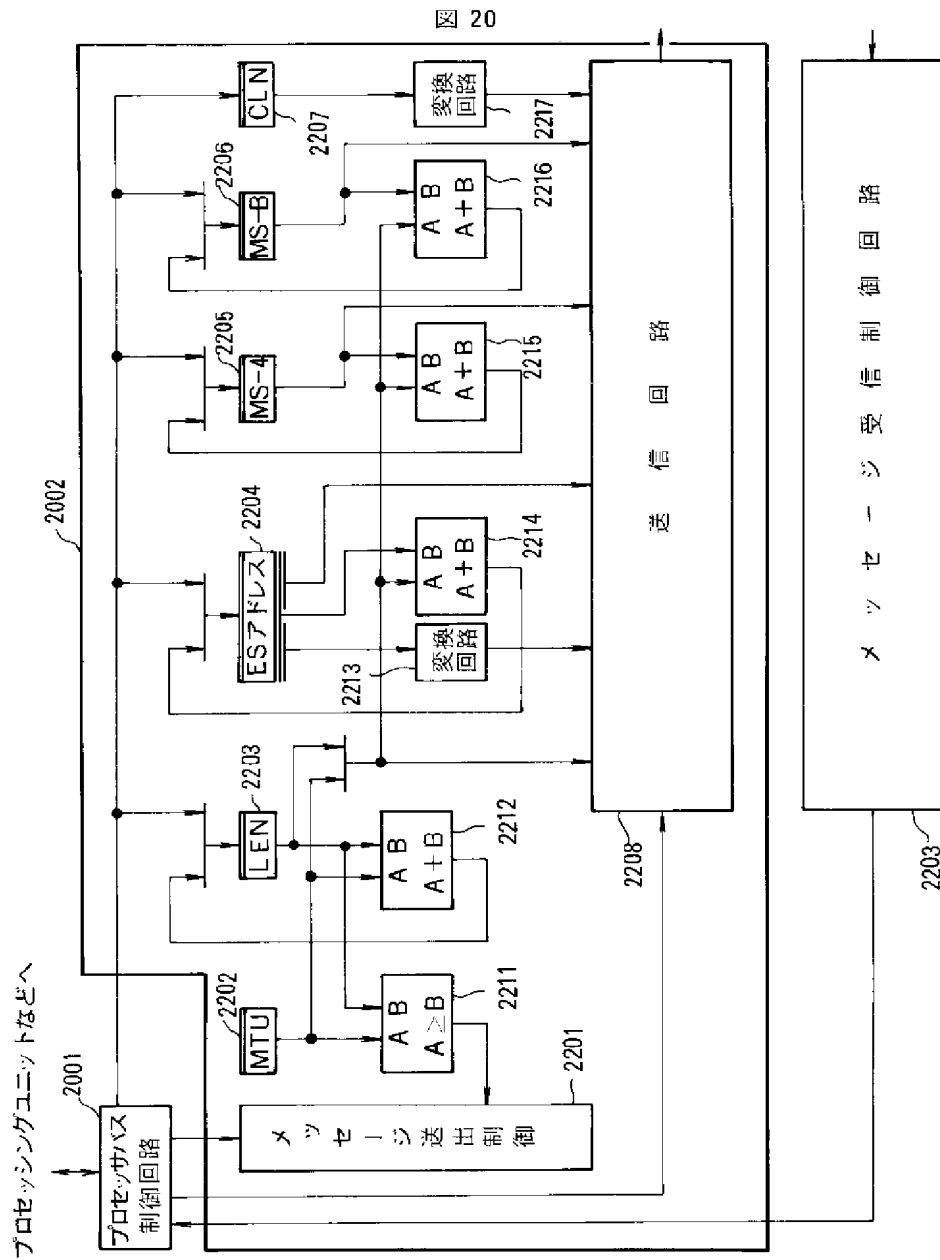
1803:ES クラスタと乗り換えスイッチ

【図19】

図 19



【図20】



フロントページの続き

(72)発明者 濱中 直樹  
東京都国分寺市東恋ヶ窪1丁目280番地  
株式会社日立製作所中央研究所内

(72)発明者 千葉 寛之  
東京都国分寺市東恋ヶ窪1丁目280番地  
株式会社日立製作所中央研究所内

(72)発明者 樋口 達雄  
東京都国分寺市東恋ヶ窪1丁目280番地  
株式会社日立製作所中央研究所内

(72)発明者 首藤 信一  
東京都国分寺市東恋ヶ窪1丁目280番地  
株式会社日立製作所中央研究所内

(72)発明者 緒方 康洋  
東京都小平市上水本町5丁目20番1号 日  
立超エル・エス・アイ・エンジニアリング  
株式会社内

(72)発明者 武内 茂雄  
東京都小平市上水本町5丁目20番1号 日  
立超エル・エス・アイ・エンジニアリング  
株式会社内

(72)発明者 鳥羽 達  
東京都小平市上水本町5丁目20番1号 日  
立超エル・エス・アイ・エンジニアリング  
株式会社内